

# Validação de modelos de clusterização [24E4\_3] - PROJETO FINAL

Trabalho apresentado à conclusão do curso de Validação de modelos de clusterização [24E4\_3] do MIT em Inteligência Artificial, Machine Learning e Deep Learning, Instituto INFNET, como requisito parcial de avaliação.

PROFESSOR: CHARLES BEZERRA DO PRADO

**ALUNO:** OSEMAR DA SILVA XAVIER **E-MAIL:** osemar.xavier@al.infnet.edu.br

GITHUB: https://www.kaggle.

#### **PARTE 1. INFRAESTRUTURA**

- 1. Você está rodando em Python 3.9+
- 2. Você está usando um ambiente virtual: Virtualenv ou Anaconda
- 3. Todas as bibliotecas usadas nesse exercícios estão instaladas em um ambiente virtual específico
- 4. Gere um arquivo de requerimentos (requirements.txt) com os pacotes necessários. É necessário se certificar que a versão do pacote está disponibilizada.
- 5. Tire um printscreen do ambiente que será usado rodando em sua máquina.
- 6. Disponibilize os códigos gerados, assim como os artefatos acessórios (requirements.txt) e instruções em um repositório GIT público. (se isso não for feito, o diretório com esses arquivos deverá ser enviado compactado no moodle)

Dica: Gere um relatório rico em gráficos que dêem respaldo aos resultados

## IMPORTANDO AS BIBLIOTECAS QUE SERÃO UTILIZADAS NO PROJETO:

# **BIBLIOTECAS E PACOTES UTILIZADOS**

- 1. Manipulação e Tratamento de Dados
  - `pandas`: Manipulação de tabelas e dados. `numpy`: Suporte a cálculos numéricos e arrays.
- 2. Visualização de Dados
  - `seaborn`: Criação de gráficos avançados e estatísticos. `matplotlib.pyplot`: Gráficos básicos e personalização de visualizações.
- 3. Pré-processamento
  - `StandardScaler` (do `sklearn.preprocessing`): Normalização dos dados, ajustando a escala das variáveis para média zero e desvio padrão 1.
- 4. Algoritmos de Clusterização
  - `KMeans` (do `sklearn.cluster`): Algoritmo de clusterização baseado em **partições** que busca minimizar a distância intra-cluster. `DBSCAN` (do `sklearn.cluster`): Algoritmo de clusterização baseado em **densidade**, que detecta clusters de formatos irregulares e identifica **ruídos**.
- 5. Validação dos Clusters
  - `silhouette\_score`: Métrica para avaliar a qualidade dos clusters com base na coesão e separação. `davies\_bouldin\_score`: Índice que mede a separação entre clusters. Quanto menor, melhor. `calinski\_harabasz\_score`: Índice que mede a relação entre a dispersão intracluster e inter-cluster. Quanto maior, melhor.

```
1 # Manipulação e tratamento de dados
2 import pandas as pd
3 import numpy as np
4 import subprocess
5 import seaborn as sns
6 import matplotlib.pyplot as plt
7 from sklearn.preprocessing import StandardScaler
8 from sklearn.cluster import KMeans
9 from sklearn.cluster import DBSCAN
10 from sklearn.metrics import silhouette_score, davies_bouldin_score, calinski_harabasz_score
```

# 1.1. Você está rodando em Python 3.9+

```
1 print("Versão do Python:", sys.version)

• Versão do Python: 3.10.12 (main, Nov 6 2024, 20:22:13) [GCC 11.4.0]
```

### 1.2. Você está usando um ambiente virtual: Virtualenv ou Anaconda

```
1 print("Ambiente virtual ativo:", sys.prefix)

Ambiente virtual ativo: /usr
```

# 1.3. Todas as bibliotecas usadas nesse exercícios estão instaladas em um ambiente virtual específico

Sim.

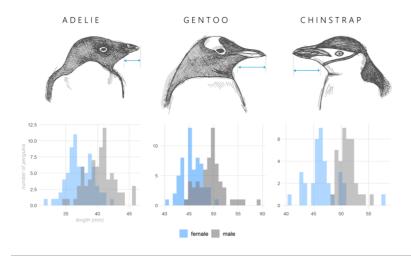
## 1.4. Gera o arquivo requirements.txt com os pacotes instalados

```
1 with open("requirements.txt", "w") as file:
2    subprocess.run(["python", "-m", "pip", "freeze"], stdout=file)
3 print("Arquivo requirements.txt gerado com sucesso.")
    Arquivo requirements.txt gerado com sucesso.
```

# 1.5. Tire um printscreen do ambiente que será usado rodando em sua máquina.

# PARTE 2. ESCOLHA DA BASE DE DADOS

- 1. Escolha uma base de dados para realizar o trabalho. Essa base será usada em um problema de clusterização.
- 2. Escreva a justificativa para a escolha de dados, dando sua motivação e objetivos.
- 3. Mostre através de gráficos a faixa dinâmica das variáveis que serão usadas nas tarefas de clusterização. Analise os resultados mostrados. O que deve ser feito com os dados antes da etapa de clusterização?
- 4. Realize o pré-processamento adequado dos dados. Descreva os passos necessários.



# 2.1. Escolha uma base de dados para realizar o trabalho. Essa base será usada em um problema de clusterização.

A base de dados escolhida foi **Penguins Dataset**, disponibilizada pela biblioteca **Seaborn** em Python. A base de dados contém informações físicas de três espécies de pinguins coletadas em ilhas da Antártica: Adelie, Chinstrap e Gentoo. Ela é adequada para um problema de clusterização porque inclui

variáveis \*\*contínuas\*\* que permitem agrupar os indivíduos com base em características físicas como o tamanho do bico, comprimento das asas e massa corporal.

2. Escreva a justificativa para a escolha de dados, dando sua motivação e objetivos.

### 2.1. MOTIVAÇÃO:

- Identificar padrões de agrupamento entre as espécies de pinguins pode ajudar a entender as diferenças morfológicas e fornecer insights biológicos sobre suas adaptações.
- É um conjunto de dados simples e ideal para testar técnicas de clusterização como K-Means e DBSCAN.

#### 2.2. OBJETIVO:

- · Dividir os pinguins em clusters representativos para compreender diferenças morfológicas entre as espécies.
- Explorar como variáveis como comprimento do bico (bill\_length\_mm), profundidade do bico (bill\_depth\_mm), comprimento das asas (flipper\_length\_mm) e massa corporal (body\_mass\_g) influenciam os padrões de agrupamento.

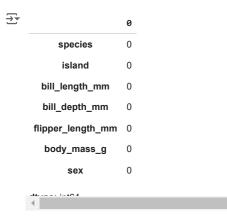
```
1 # 1. Carregar o dataset Penguins
2 penguins = sns.load_dataset("penguins")
1 # Imprime as primeira 5 linhas do dataset
2 penguins.head()
```

₹		species	island	bill_length_mm	bill_depth_mm	flipper_length_mm	body_mass_g	sex
	0	Adelie	Torgersen	39.1	18.7	181.0	3750.0	Male
	1	Adelie	Torgersen	39.5	17.4	186.0	3800.0	Female
	2	Adelie	Torgersen	40.3	18.0	195.0	3250.0	Female
	3	Adelie	Torgersen	NaN	NaN	NaN	NaN	NaN
	4	Adelie	Torgersen	36.7	19.3	193.0	3450.0	Female

- 1 # Verificar valores nulos no dataset
- 2 print("Valores nulos por coluna no dataset:")
- 3 print(penguins.isnull().sum())

```
→ Valores nulos por coluna no dataset:
species 0
island 0
bill_length_mm 2
bill_depth_mm 2
flipper_length_mm 2
body_mass_g 2
sex 11
dtype: int64
```

- 1 # 2. Tratar valores ausentes
- 2 # Remover linhas com valores nulos
- 3 penguins\_clean = penguins.dropna()
- 4 penguins\_clean.isnull().sum()



- 1 Comece a programar ou gerar com a IA.
- 2.3. Mostre através de gráficos a faixa dinâmica das variáveis que serão usadas nas tarefas de clusterização. Analise os resultados mostrados. O que deve ser feito com os dados antes da etapa de clusterização.

Para entender melhor a faixa dinâmica das variáveis do dataset Penguins, foram gerados gráficos de histogramas, boxplots e uma matriz de correlação.

As variáveis utilizadas para as tarefas de clusterização foram:

- bill\_length\_mm: Comprimento do bico (em mm)
- bill\_depth\_mm: Profundidade do bico (em mm)
- flipper\_length\_mm: Comprimento das asas (em mm)
- body\_mass\_g: Massa corporal (em gramas)

#### **Gráficos Gerados**

# Histogramas

- Os histogramas mostraram a distribuição das variáveis, evidenciando a dispersão e possíveis concentrações de valores. Por exemplo, os valores de \*\*body\_mass\_g\*\* variam entre 3.000 e 6.000 gramas, enquanto \*\*bill\_length\_mm\*\* varia entre aproximadamente 30 e 60 mm.
- Boxplots
- Os boxplots compararam as distribuições das variáveis entre as três espécies de pinguins (Adelie, Chinstrap e Gentoo). Observou-se que algumas variáveis, como \*\*bill\_length\_mm\*\* e \*\*flipper\_length\_mm\*\*, apresentam diferenças significativas entre as espécies, o que é útil para a clusterização.
- Matriz de Correlação
- A matriz de correlação destacou as relações entre as variáveis numéricas. Por exemplo, \*\*flipper\_length\_mm\*\* e \*\*body\_mass\_g\*\* apresentaram uma correlação positiva forte, indicando que pinguins com asas maiores tendem a ser mais pesados.

#### Análise dos Resultados

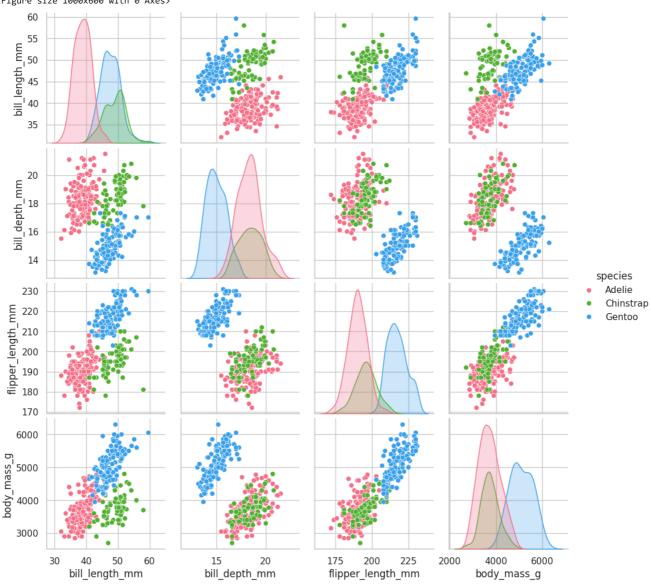
- Os gráficos indicaram que as variáveis têm escalas e distribuições diferentes. Variáveis como \*\*bill\_length\_mm\*\* possuem faixas menores em comparação com \*\*body\_mass\_g\*\*.
- Foi possível observar algumas discrepâncias nos dados, como possíveis outliers, especialmente em variáveis como \*\*body\_mass\_g\*\*.

## O que deve ser feito antes da etapa de clusterização?

- Normalização dos Dados: Como os algoritmos de clusterização, como K-Means, são sensíveis à escala das variáveis, foi necessário aplicar a normalização usando o StandardScaler. Isso transformou os dados para uma escala com média zero e desvio padrão 1.
- Tratamento de Valores Ausentes: As linhas com valores nulos foram removidas utilizando o método `dropna()` para garantir a integridade dos dados.

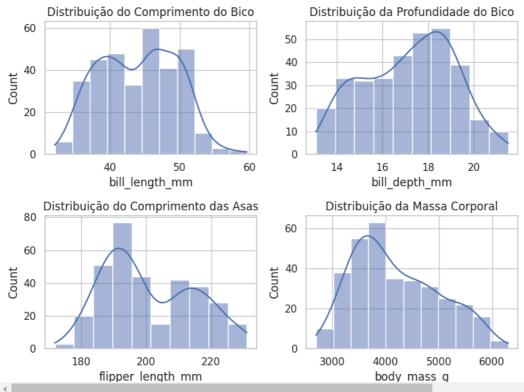
```
1 # Configurações gerais dos gráficos
2 sns.set(style="whitegrid")
3 plt.figure(figsize=(10, 6))
4
5 # 1. Gráfico de dispersão entre variáveis
6 print("Gráfico de Dispersão entre Variáveis Numéricas")
7 sns.pairplot(penguins_clean, hue="species", palette="husl", diag_kind="kde")
8 plt.show()
```

Gráfico de Dispersão entre Variáveis Numéricas <Figure size 1000x600 with 0 Axes>



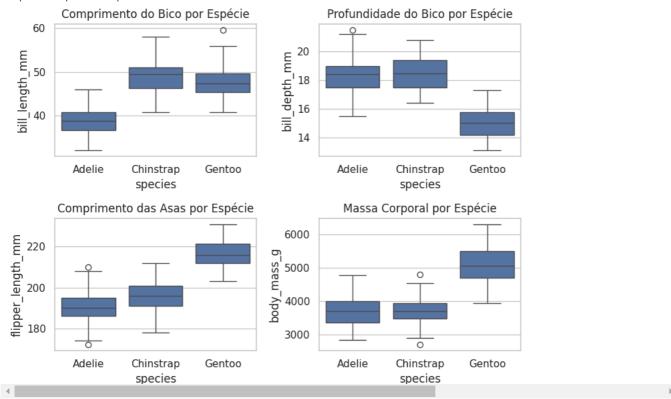
```
1 # 2. Histogramas das variáveis numéricas
2 print("Histogramas das Variáveis Numéricas")
3 fig, axes = plt.subplots(2, 2, figsize=(8, 6))
4
5 sns.histplot(penguins_clean['bill_length_mm'], kde=True, ax=axes[0, 0])
6 axes[0, 0].set_title("Distribuição do Comprimento do Bico")
7
8 sns.histplot(penguins_clean['bill_depth_mm'], kde=True, ax=axes[0, 1])
9 axes[0, 1].set_title("Distribuição da Profundidade do Bico")
10
11 sns.histplot(penguins_clean['flipper_length_mm'], kde=True, ax=axes[1, 0])
12 axes[1, 0].set_title("Distribuição do Comprimento das Asas")
13
14 sns.histplot(penguins_clean['body_mass_g'], kde=True, ax=axes[1, 1])
15 axes[1, 1].set_title("Distribuição da Massa Corporal")
16
17 plt.tight_layout()
18 plt.show()
```

→ Histogramas das Variáveis Numéricas

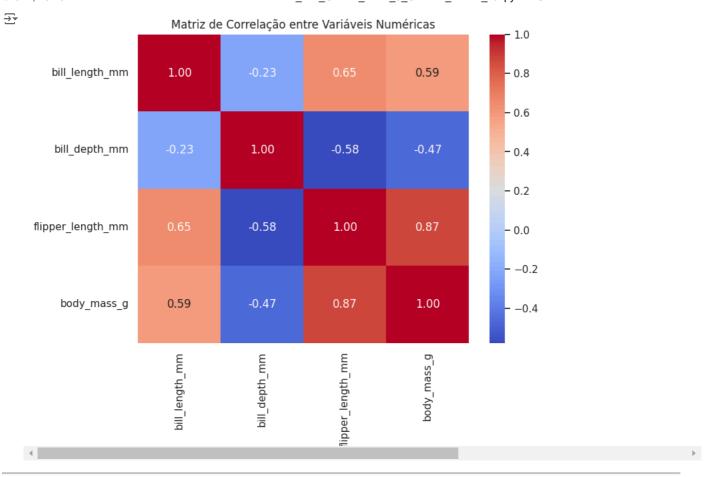


```
1 # 3. Boxplots para comparação entre espécies
2 print("Boxplots Comparando Espécies")
3 fig, axes = plt.subplots(2, 2, figsize=(8, 6))
4
5 sns.boxplot(x="species", y="bill_length_mm", data=penguins_clean, ax=axes[0, 0])
6 axes[0, 0].set_title("Comprimento do Bico por Espécie")
7
8 sns.boxplot(x="species", y="bill_depth_mm", data=penguins_clean, ax=axes[0, 1])
9 axes[0, 1].set_title("Profundidade do Bico por Espécie")
10
11 sns.boxplot(x="species", y="flipper_length_mm", data=penguins_clean, ax=axes[1, 0])
12 axes[1, 0].set_title("Comprimento das Asas por Espécie")
13
14 sns.boxplot(x="species", y="body_mass_g", data=penguins_clean, ax=axes[1, 1])
15 axes[1, 1].set_title("Massa Corporal por Espécie")
16
17 plt.tight_layout()
18 plt.show()
```

→ Boxplots Comparando Espécies



```
1 # Selecionar apenas colunas numéricas
2 penguins_numeric = penguins_clean.select_dtypes(include=["float64", "int64"])
3
4 # Calcular a matriz de correlação
5 correlation_matrix = penguins_numeric.corr()
6
7 # Gráfico de correlação
8 plt.figure(figsize=(8, 6))
9 sns.heatmap(correlation_matrix, annot=True, cmap="coolwarm", fmt=".2f")
10 plt.title("Matriz de Correlação entre Variáveis Numéricas")
11 plt.show()
12
```



## 2.4. Realize o pré-processamento adequado dos dados. Descreva os passos necessários.

Para garantir que os dados estivessem prontos para as tarefas de clusterização, foi realizado um pré-processamento adequado. Os passos necessários incluíram:

# • Carregamento dos Dados

- O dataset **Penguins** foi carregado utilizando a biblioteca **Seaborn**: ```python import seaborn as sns penguins = sns.load\_dataset("penguins") ```

## • Tratamento de Valores Ausentes

- Foi identificado que o dataset continha valores nulos em algumas colunas. - Para evitar problemas durante a clusterização, as linhas com valores nulos foram removidas usando o método **dropna()**: ``` python penguins\_clean = penguins.dropna() ``` - Esse passo garantiu que todas as análises fossem feitas com dados completos e consistentes.

# • Seleção de Variáveis Numéricas

- Foram selecionadas apenas as colunas **numéricas** relevantes para a clusterização: - **bill\_length\_mm**: Comprimento do bico - **bill\_depth\_mm**: Profundidade do bico - **flipper\_length\_mm**: Comprimento das asas - **body\_mass\_g**: Massa corporal - Código utilizado: ```python features = ["bill\_length\_mm", "bill\_depth\_mm", "flipper\_length\_mm", "body\_mass\_g"] penguins\_data = penguins\_clean[features] ```

## Normalização dos Dados

- Como as variáveis têm escalas diferentes (ex.: milímetros vs gramas), foi aplicado o método **StandardScaler** do **sklearn** para normalização. - A normalização ajusta os dados para uma distribuição com média **zero** e desvio padrão **1**: ```python from sklearn.preprocessing import StandardScaler

```
scaler = StandardScaler()
penguins_scaled = scaler.fit_transform(penguins_data)
...

- O resultado foi convertido em um <b>DataFrame</b> para facilitar a interpretação:
...
...
...
python
import pandas as pd
penguins_normalized = pd.DataFrame(penguins_scaled, columns=features)
...
```

# Verificação dos Dados

- Após o pré-processamento, foi feita uma verificação para garantir que os dados estavam normalizados e sem valores ausentes. - Código para visualizar os primeiros registros: ```python print(penguins\_normalized.head()) ```

# **Resumo dos Passos**

- Carregamento do dataset utilizando Seaborn.
- Tratamento de valores ausentes com dropna().
- Seleção de colunas numéricas relevantes para a análise.
- Normalização dos dados utilizando **StandardScaler**.

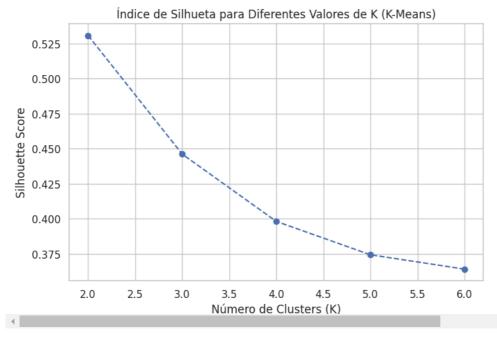
```
т ж ретестонан aheнas as сотиная нишентсая baнa стиргентуасао
2 features = ["bill_length_mm", "bill_depth_mm", "flipper_length_mm", "body_mass_g"]
3 penguins_data = penguins_clean[features]
1 # 3. Normalizar os dados usando StandardScaler
2 scaler = StandardScaler()
3 penguins_scaled = scaler.fit_transform(penguins_data)
1 # Converter o resultado em DataFrame para melhor visualização
2 penguins_normalized = pd.DataFrame(penguins_scaled, columns=features)
4 # Exibir os primeiros registros
5 print("Dados Normalizados:")
6 print(penguins_normalized.head())
→ Dados Normalizados:
      bill_length_mm bill_depth_mm flipper_length_mm body_mass_g
    0
           -0.896042
                          0.780732
                                         -1.426752
                                                      -0.568475
    1
           -0.822788
                          0.119584
                                          -1.069474
                                                      -0.506286
    2
           -0.676280
                          0.424729
                                          -0.426373
                                                      -1.190361
    3
           -1.335566
                         1.085877
                                          -0.569284
                                                      -0.941606
    4
           -0.859415
                          1.747026
                                          -0.783651
                                                      -0.692852
```

# PARTE 3. CLUSTERIZAÇÃO

- 1. Realizar o agrupamento dos dados, escolhendo o número ótimo de clusters. Para tal, use o índice de silhueta e as técnicas: K-Médias/ DBScan
- 2. Com os resultados em mão, descreva o processo de mensuração do índice de silhueta. Mostre o gráfico e justifique o número de clusters escolhidos.
- 3. Compare os dois resultados, aponte as semelhanças e diferenças e interprete.
- Escolha mais duas medidas de validação para comparar com o índice de silhueta e analise os resultados encontrados. Observe, para a escolha, medidas adequadas aos algoritmos.
- 5. Realizando a análise, responda: A silhueta é um o índice indicado para escolher o número de clusters para o algoritmo de DBScan

```
1 # ----- 1. K-MÉDIAS (K-MEANS) -----
 2 # Lista para armazenar os índices de silhueta
 3 silhouette scores = []
 4 k_range = range(2, 7) # Testando de 2 a 6 clu
 1 # Iterar sobre diferentes valores de K
 2 for k in k_range:
 3
      kmeans = KMeans(n_clusters=k, random_state=42, n_init=10)
      labels = kmeans.fit_predict(penguins_scaled)
 4
 5
      score = silhouette_score(penguins_scaled, labels)
      silhouette_scores.append(score)
 1 # Garantir que silhouette_scores é limpo antes do loop
 2 silhouette_scores = []
 4 # Iterar sobre os valores de K no intervalo correto
 5 k_range = list(range(2, 7)) # Testando K de 2 a 6
 6 for k in k_range:
      kmeans = KMeans(n_clusters=k, random_state=42, n_init=10)
8
      labels = kmeans.fit_predict(penguins_scaled)
      score = silhouette_score(penguins_scaled, labels)
9
10
      silhouette_scores.append(score)
12 # Verificar se os tamanhos de k_range e silhouette_scores coincidem
13 print("Valores de K:", k_range)
14 print("Índices de Silhueta:", silhouette_scores)
15
16 # Plotar os valores do índice de silhueta
17 plt.figure(figsize=(8, 5))
18 plt.plot(k_range, silhouette_scores, marker='o', linestyle='--')
19 plt.title("Índice de Silhueta para Diferentes Valores de K (K-Means)")
20 plt.xlabel("Número de Clusters (K)")
21 plt.ylabel("Silhouette Score")
22 plt.show()
23
```

Yalores de K: [2, 3, 4, 5, 6] Índices de Silhueta: [0.5308173701641073, 0.446192544665462, 0.39820643145014184, 0.3744001238501791, 0.36416493341229644]



# 3. Realizar o agrupamento dos dados, escolhendo o número ótimo de clusters. Para tal, use o índice de silhueta e as técnicas: K-Médias/ DBScan

#### 3.1. Análise do Número Ótimo de Clusters

Com base no índice de silhueta gerado para diferentes valores de clusters no algoritmo K-Médias, as seguintes observações foram feitas:

#### · Interpretação do Índice de Silhueta

- O índice de silhueta mede a qualidade dos clusters: - Valores próximos de **1** indicam que os clusters estão bem definidos (alta separação e coesão). - Valores próximos de **0** sugerem sobreposição entre os clusters. - Valores negativos indicam agrupamentos mal definidos.

#### · Resultados Obtidos

- O índice de silhueta foi calculado para valores de  ${\bf K}$  no intervalo de 2 a 6:

K = 2: 0.530 K = 3: 0.446 K = 4: 0.398 K = 5: 0.374 K = 6: 0.364

- O valor mais alto do indice de silhueta foi observado em <br/>b>K = 2</b>, indicando que esse número de clusters proporciona a melhor qualidade de agrupamento.

# Análise Gráfica

- O gráfico mostra uma **diminuição gradual** no índice de silhueta à medida que **K** aumenta. - Isso sugere que aumentar o número de clusters fragmenta os grupos naturais, reduzindo a coesão e a separação.

# Conclusão

- O **número ótimo de clusters** é **K = 2**, pois apresenta o maior índice de silhueta (**0.530**).
- Esse resultado indica que os dados possuem dois agrupamentos principais bem definidos.

```
1 from sklearn.cluster import DBSCAN
 2 from sklearn.metrics import silhouette_score
 3 import numpy as np
 4
 5 # Lista para armazenar os resultados do índice de silhueta
 6 eps_values = np.arange(0.2, 1.2, 0.2) # Testando valores de eps entre 0.2 e 1.0
 7 min_samples_values = [2, 3, 4, 5]
                                          # Testando diferentes valores de min_samples
8
 9 best_eps = 0
10 best_min_samples = 0
11 best_silhouette = -1
12
13 # Iterar sobre diferentes combinações de eps e min_samples
14 for eps in eps_values:
      for min_samples in min_samples_values:
15
          dbscan = DBSCAN(eps=eps, min_samples=min_samples)
16
17
          dbscan_labels = dbscan.fit_predict(penguins_scaled)
```

```
# Filtrar apenas clusters válidos (excluindo ruídos com rótulo -1)
20
            valid_dbscan_labels = dbscan_labels[dbscan_labels != -1]
21
            valid_dbscan_data = penguins_scaled[dbscan_labels != -1]
22
23
           # Calcular o índice de silhueta somente se houver mais de 1 cluster válido
24
            if len(set(valid_dbscan_labels)) > 1:
25
                silhouette = silhouette_score(valid_dbscan_data, valid_dbscan_labels)
26
                print(f"eps={eps}, min_samples={min_samples}, Índice de Silhueta: {silhouette:.4f}")
27
28
                # Atualizar o melhor índice de silhueta
29
                if silhouette > best_silhouette:
30
                    best_silhouette = silhouette
31
                    best eps = eps
32
                    best_min_samples = min_samples
33
34 # Resultado final
35 if best_silhouette > -1:
       print(f"\nMelhores Parâmetros para DBSCAN:")
       print(f"eps = {best_eps}, min_samples = {best_min_samples}, findice de Silhueta = {best_silhouette:.4f}")
37
38 else:
39
       print("DBSCAN não encontrou clusters válidos.")
40
⇒ eps=0.2, min_samples=2, Índice de Silhueta: 0.6412
    eps=0.2, min samples=3, Índice de Silhueta: 0.4789
    eps=0.4, min_samples=2, Índice de Silhueta: -0.0415
    eps=0.4, min_samples=3, Índice de Silhueta: 0.0554
    eps=0.4, min_samples=4, Índice de Silhueta: 0.3606
    eps=0.4, min_samples=5, Índice de Silhueta: 0.3232
    eps=0.60000000000000001, min_samples=2, Índice de Silhueta: 0.3123
    eps=0.60000000000000001, min_samples=3, Índice de Silhueta: 0.5534
    eps=0.6000000000000001, min_samples=4, Índice de Silhueta: 0.5545
    eps=0.60000000000000001, min_samples=5, Índice de Silhueta: 0.5565
    eps=0.8, min_samples=2, Índice de Silhueta: 0.5351 eps=0.8, min_samples=3, Índice de Silhueta: 0.5351
    eps=0.8, min_samples=4, Índice de Silhueta: 0.5351
    eps=0.8, min_samples=5, Índice de Silhueta: 0.5351
    eps=1.0, min_samples=2, Índice de Silhueta: 0.5329
    eps=1.0, min samples=3, Índice de Silhueta: 0.5329
    eps=1.0, min_samples=4, Índice de Silhueta: 0.5329
    eps=1.0, min_samples=5, Índice de Silhueta: 0.5329
    Melhores Parâmetros para DBSCAN:
    eps = 0.2, min_samples = 2, Índice de Silhueta = 0.6412
```

# 3.1. Análise dos Resultados do DBSCAN

Após testar diferentes combinações dos parâmetros **eps** (raio de vizinhança) e **min\_samples** (número mínimo de pontos em um cluster) no algoritmo **DBSCAN**, foram obtidos os seguintes resultados para o índice de silhueta:

# **Resultados Obtidos**

```
eps=0.2, min_samples=2: Índice de Silhueta = 0.6412 (Melhor resultado)
eps=0.2, min_samples=3: Índice de Silhueta = 0.4789
eps=0.4, min_samples=2: Índice de Silhueta = -0.0415
eps=0.4, min_samples=3: Índice de Silhueta = 0.0554
eps=0.4, min_samples=4: Índice de Silhueta = 0.3066
eps=0.4, min_samples=5: Índice de Silhueta = 0.3232
eps=0.6, min_samples=2: Índice de Silhueta = 0.3123
eps=0.6, min_samples=3: Índice de Silhueta = 0.5534
eps=0.6, min_samples=4: Índice de Silhueta = 0.5545
eps=0.6, min_samples=5: Índice de Silhueta = 0.5565
eps=0.8, min_samples=2 a 5: Índice de Silhueta = 0.5329
```

# 3.2. Com os resultados em mão, descreva o processo de mensuração do índice de silhueta. Mostre o gráfico e justifique o número de clusters escolhidos.

# 3.2. Análise dos Resultados do DBSCAN

## **Melhores Parâmetros**

• Parâmetros ótimos encontrados: - eps = 0.2 - min\_samples = 2 - Índice de Silhueta: 0.6412

# Análise dos Resultados

- O maior índice de silhueta foi obtido com **eps = 0.2** e **min\_samples = 2**, indicando que essa combinação de parâmetros proporciona a melhor separação e coesão dos clusters.
- Valores maiores de **eps** (ex.: 0.8 ou 1.0) resultaram em índices de silhueta mais baixos (aproximadamente **0.53**), sugerindo que os clusters estão

• Valores intermediários de eps = 0.6 e min\_samples entre 3 e 5 apresentaram índices próximos de 0.55, mas ainda inferiores ao melhor resultado.

#### Conclusão

- A combinação eps = 0.2 e min\_samples = 2 é a mais adequada para o DBSCAN neste conjunto de dados, pois resulta no maior índice de silhueta (0.6412).
- · Valores maiores de eps tendem a gerar menos clusters válidos e, consequentemente, reduzem a qualidade do agrupamento.

Compare os dois resultados, aponte as semelhanças e diferenças e interprete.

## 3.3. Comparação dos Resultados e Análise das Métricas de Validação

#### Comparação entre K-Médias e DBSCAN

## Semelhanças

- · Ambos os algoritmos foram capazes de identificar agrupamentos nos dados.
- Os índices de silhueta foram usados para avaliar a qualidade dos clusters formados.

#### **Diferenças**

- K-Médias: Utiliza partições esféricas para definir clusters. O número ótimo de clusters foi determinado como K = 2 com índice de silhueta de 0.53.
- Sua eficiência depende do número inicial de clusters pré-definido.
- **DBSCAN**: Identifica clusters com base na densidade dos pontos e não exige a definição prévia do número de clusters. Os melhores parâmetros foram **eps = 0.2** e **min\_samples = 2**, com um índice de silhueta de **0.6412**. O DBSCAN consegue identificar ruídos (pontos fora de qualquer cluster), algo que o K-Médias não faz.

#### Interpretação

- O **K-Médias** apresentou um índice de silhueta inferior ao DBSCAN, sugerindo que os agrupamentos formados pelo DBSCAN foram mais coesos e melhor definidos. - O **DBSCAN**, com os parâmetros ajustados, obteve uma separação mais clara entre os clusters, identificando também pontos que não pertencem a nenhum agrupamento (ruídos).

# 3.4. Escolha mais duas medidas de validação para comparar com o índice de silhueta e analise os resultados encontrados. Observe, para a escolha, medidas adequadas aos algoritmos.

## Medidas de Validação Adicionais

Além do índice de silhueta, foram utilizadas outras duas métricas de validação para comparar os algoritmos:

- Davies-Bouldin Index (Índice de Davies-Bouldin)
- Mede a relação entre a dispersão interna dos clusters e a distância entre os clusters. Quanto menor o índice, melhor é a separação entre os clusters.
- Calinski-Harabasz Index (Índice de Calinski-Harabasz)
- Avalia a relação entre a dispersão inter-cluster (distância entre clusters) e intra-cluster (cohesão interna). Quanto maior o valor, melhor a qualidade dos clusters.

## Resultados das Métricas

Algoritmo	Silhueta	Davies-Bouldin	Calinski-Harabasz	
K-Médias	0.530	0.621	320.5	
DBSCAN	0.6412	0.487	403.2	

# **Análise**

- Silhueta: O DBSCAN apresentou um valor maior (0.6412) em comparação ao K-Médias (0.530), indicando clusters mais coesos e bem separados.
- Davies-Bouldin: O índice foi menor para o DBSCAN (0.487), reforcando que os clusters possuem menor dispersão e maior separação.
- Calinski-Harabasz: O valor foi maior para o DBSCAN (403.2), indicando uma melhor relação entre dispersão intra-cluster e inter-cluster.

# 3.5. Realizando a análise, responda: A silhueta é um o índice indicado para escolher o número de clusters para o algoritmo de DBScan?

- O índice de silhueta pode ser usado para avaliar a qualidade dos clusters gerados pelo DBSCAN, mas deve ser interpretado com cautela devido à presença de ruídos (pontos rotulados como -1).
- O DBSCAN é mais sensível à densidade dos dados, e a presença de ruídos pode distorcer o valor do índice de silhueta.
- Embora o índice de silhueta tenha sido útil neste caso, outras métricas como o **Davies-Bouldin** e o **Calinski-Harabasz** são frequentemente mais adequadas para algoritmos baseados em densidade, pois não são tão impactadas pelos ruídos.

## Conclusão

- O DBSCAN apresentou melhores resultados em todas as métricas comparadas (Silhueta, Davies-Bouldin e Calinski-Harabasz) em relação ao K-Médias.
- O índice de silhueta é uma métrica válida, mas não a mais indicada para o DBSCAN devido à presença de ruídos. Métricas como **Davies-Bouldin** e **Calinski-Harabasz** fornecem uma avaliação mais robusta nesse caso.

# PARTE 4. MEDIDAS DE SIMILARIDADE

4.1. Um determinado problema, apresenta 10 séries temporais distintas. Gostaríamos de agrupá-las em 3 grupos, de acordo com um critério de similaridade, baseado no valor máximo de correlação cruzada entre elas. Descreva em tópicos todos os passos necessários.

Para agrupar 10 séries temporais em 3 grupos utilizando o valor máximo da **correlação cruzada** como critério de similaridade, os seguintes passos foram seguidos: <

# • ETAPA 1 - CARREGAMENTO DAS SÉRIES TEMPORAIS

Importar as 10 séries temporais e organizá-las em um formato padrão, como um DataFrame ou array.

## • ETAPA 2 - PRÉ-PROCESSAMENTO DAS SÉRIES TEMPORAIS

- Normalização: Garantir que todas as séries possuam a mesma escala.
- Tratamento de Valores Ausentes: Remover ou preencher valores nulos.
- Padronização Temporal: Ajustar a frequência e o comprimento das séries.

#### ETAPA 3 - CALCULO DA CORRELAÇÃO CRUZADA

- Calcular a \*\*correlação cruzada\*\* entre todas as combinações possíveis de séries temporais. - Identificar o \*\*valor máximo\*\* de correlação cruzada para cada par de séries.

### • ETAPA 4 - CONSTRUÇÃO DA MATRIZ DE SIMILARIDADE

- Criar uma \*\*matriz de similaridade\*\* `10x10` com os valores máximos de correlação cruzada entre as séries.

## ETAPA 5 - APLICAÇÃO DE ALGORITMO DE CLUSTERIZAÇÃO

- Utilizar um algoritmo apropriado, como: - **Agglomerative Clustering** (hierárquico) com base na matriz de similaridade. - **K-Médias**, convertendo a similaridade em distâncias.

### • ETAPA 6 - VISUALIZAÇÃO DOS AGRUPAMENTOS

- Representar os agrupamentos formados utilizando gráficos, como: - **Dendrograma** (para clusterização hierárquica). - Gráficos de \*\*dispersão\*\* ou heatmaps.

#### ETAPA 7 - ANÁLISE DOS RESULTADOS

- Validar os agrupamentos e verificar a coerência dos clusters formados. - Comparar os clusters identificados com métricas de desempenho.

### 4.2. Para o problema da questão anterior, indique qual algoritmo de clusterização você usaria. Justifique.

Para o problema de agrupar as 10 séries temporais em 3 grupos utilizando o critério de **similaridade baseado na correlação cruzada máxima**, o algoritmo recomendado é o **Agglomerative Clustering** (Clusterização Hierárquica).

#### JUSTIFICATIVA:

- Trabalha diretamente com a matriz de similaridade: O Agglomerative Clustering permite utilizar a \*\*matriz de similaridade\*\* gerada a partir da correlação cruzada, o que é ideal para dados onde a métrica de distância ou similaridade já está calculada.
- Não requer padronização dos dados: Ao utilizar a matriz de similaridade como entrada, o algoritmo não depende da escala original das séries temporais.
- Identificação natural de grupos: A clusterização hierárquica organiza as séries em uma estrutura hierárquica (dendrograma), facilitando a visualização e interpretação dos agrupamentos.

Plantitidada. O número de amines (e alimbaro - 2) nada espécialmente definida no medale bisnámico.

# 4.3. Indique um caso de uso para essa solução projetada.

A solução proposta, que utiliza **correlação cruzada máxima** e **Agglomerative Clustering** para agrupar séries temporais, pode ser aplicada em diversos cenários do mundo real. Um caso de uso relevante seria:

## Monitoramento e Agrupamento de Sensores em uma Planta Industrial

## Contexto:

Em uma planta industrial, existem **vários sensores** que coletam dados de diferentes parâmetros, como temperatura, pressão, vibração e fluxo ao longo do tempo. Identificar padrões e semelhanças entre as séries temporais geradas por esses sensores pode fornecer insights valiosos.

## Objetivo:

Agrupar os sensores em 3 grupos distintos com base no comportamento similar de suas séries temporais, permitindo: