

**UJIAN AKHIR SEMESTER:  
CLUSTERING RENTANG SUHU CUACA JAKARTA TAHUN  
2021-2023 MENGGUNAKAN ALGORITMA K-MEAN DAN  
GAUSSIAN MIXTURE MODELS**



**UNIVERSITAS  
BUDI LUHUR**

Disusun oleh :

Muchamad Angga Dwi Wahyu      2112501339

**UNIVERSITAS BUDI LUHUR  
FAKULTAS TEKNOLOGI INFORMASI  
JAKARTA  
2022/2023**

## DAFTAR ISI

DAFTAR ISI.....	i
BAB I PENDAHULUAN.....	1
1.1. Latar Belakang.....	1
1.2. Perumusan Masalah.....	1
1.3. Tujuan.....	1
BAB II METODE PENELITIAN.....	2
2.1. Metode Yang Diganakan.....	2
BAB III INFORMASI DATASET.....	3
3.1. Sample Data.....	3
3.2. Meta Data.....	4
3.3. Transformasi Data.....	4
BAB IV PERSIAPAN DATA.....	5
4.1. Data Cleaning.....	5
BAB V VISUALISASI DATA.....	6
5.1. Grafik Time Series.....	6
BAB VI PEMROSESAN DATA.....	7
6.1. Pemrosesan K-Means.....	7
6.2. Pemrosesan Gaussian Mixture Models.....	8
BAB VII HASIL.....	9
7.1 Hasil Metode K-Means.....	9
7.2 Hasil Metode Gaussian Mixture Models.....	10
7.3 Perbandingan pemetaan data hasil metode K-Means dengan GMM.....	15
BAB VIII PENUTUP.....	12
8.1 Kesimpulan.....	10

## **BAB I PENDAHULUAN**

### **1.1. Latar Belakang**

Cuaca memainkan peran kunci dalam kehidupan sehari-hari dan dapat memengaruhi berbagai aspek kehidupan manusia, termasuk pertanian, transportasi, dan kesehatan. Pemahaman mendalam tentang pola cuaca sangat diperlukan untuk mengantisipasi dan mengelola dampaknya.

Jakarta, sebagai ibu kota Indonesia, mengalami variasi cuaca yang signifikan sepanjang tahun. Dari musim hujan hingga musim kemarau, pemahaman tentang perubahan cuaca dalam rentang waktu tertentu dapat memberikan wawasan tentang tren iklim dan perubahan yang mungkin terjadi.

Dalam menghadapi variasi cuaca yang mengalami perubahan yang signifikan, Penelitian ini memiliki dua tujuan utama: pertama, membentuk kelompok-kelompok suhu cuaca harian yang serupa menggunakan algoritma K-Means dan Gaussian Mixture Models (GMM); kedua, membandingkan kinerja kedua algoritma ini dengan menggunakan metrik evaluasi Silhouette Score. Hasil penelitian diharapkan dapat memberikan wawasan yang lebih baik tentang pola kondisi suhu cuaca di Jakarta serta menentukan algoritma yang paling optimal untuk tugas pengelompokan suhu cuaca harian.

### **1.2. Perumusan Masalah**

Berdasarkan uraian latar belakang maka diuraikan perumusan masalah sebagai berikut :

1. Bagaimana membagi data suhu cuaca harian berdasarkan nilai suhu minimal dan maksimal ke dalam kelompok-kelompok yang serupa?
2. Manakah algoritma yang paling optimal untuk digunakan sebagai pembuat cluster untuk menyelesaikan permasalahan 1?

### **1.3. Tujuan**

Berdasarkan uraian perumusan masalah maka diuraikan tujuan dari artikel ilmiah ini sebagai berikut :

1. Membuat kluster rentang kondisi suhu cuaca di Jakarta berdasarkan nilai suhu cuaca minimal dan suhu cuaca maksimal Jakarta tahun 2021-2023.
2. Mengetahui algoritma yang optimal antar K-Mean dan GMM untuk membuat cluster rentang suhu cuaca Jakarta tahun 2021-2023 berdasarkan Silhouette Score.

## **BAB II**

### **METODE PENELITIAN**

#### **2.1. Metode Yang Digunakan**

Dalam penelitian ini akan digunakan dua metode machine learning yaitu metode K-Means dan Gaussian Mixture Models:

1. K-Means : Merupakan metode yang menggunakan pendekatan partisi untuk membentuk kelompok data. Cocok untuk data dengan kelompok yang jelas dan sederhana, memberikan solusi clustering dengan komputasi yang efisien.
2. Gaussian Mixture Models : Memiliki tujuan yang sama dengan K-Means tetapi menggunakan pendekatan probabilistik di mana setiap data dapat diberikan probabilitas untuk menjadi bagian dari setiap kelompok. Ideal untuk mengidentifikasi pola suhu cuaca yang mungkin bersifat kompleks dan tidak teratur.
3. Penilaian performa klastering, metode Silhouette akan dilibatkan untuk mengetahui algoritma yang optimal antar K-Mean dan GMM.

## BAB III

### INFORMASI DATASET

### 3.1. Sample Data Asli

[illegible]

### 3.2. Meta Data

Berikut adalah *meta data* dari dataset yang dipakai :

1. Sumber : <https://weather.visualcrossing.com/VisualCrossingWebServices/rest/services/timeline>
2. Author : visualcrossing
3. Ukuran File : 506 kb
4. Jumlah Baris : 1001
5. Jumlah Kolom : 33

### 3.3. Transformasi Data

Dari data set original yang akan digunakan hanya kolom “name”, “datetime”, “tempmax”, dan “tempmin” yang akan digunakan, disini bisa dilihat bahwa penggunaan format data dalam kolom “tempmax”, dan “tempmin” masih menggunakan titik yang merupakan format general/text, di tahap selanjutnya data cleaning akan diubah menjadi format number menggunakan koma(.). Berikut data yang sudah di transformasi :

name	datetime	tempmax	tempmin
Jakarta	07/04/2021	30.6	25
Jakarta	08/04/2021	33	25
Jakarta	09/04/2021	32	26
Jakarta	10/04/2021	32	23.7
Jakarta	11/04/2021	32.7	24.6
Jakarta	12/04/2021	33	24.3
Jakarta	13/04/2021	33	25
Jakarta	14/04/2021	32	24.3
Jakarta	15/04/2021	32	24
Jakarta	16/04/2021	32	24.3
Jakarta	17/04/2021	32	25.3
Jakarta	18/04/2021	32	25.3
Jakarta	19/04/2021	32.7	25
Jakarta	20/04/2021	32.7	24
Jakarta	21/04/2021	33.7	24.7
Jakarta	22/04/2021	33	24.5
Jakarta	23/04/2021	32.7	24.3
Jakarta	24/04/2021	32.7	24.3
Jakarta	25/04/2021	33	25.3
Jakarta	26/04/2021	34	26
Jakarta	27/04/2021	34	25.3
Jakarta	28/04/2021	32	24.3
Jakarta	29/04/2021	33	24.3

## BAB IV PERSIAPAN DATA

### 4.1. Data Cleaning (Menggunakan Python)

#### 8.1 Modul Yang Digunakan :

- Pandas

#### 8.2 Kode Program :

```
import pandas as pd

# Membaca file Excel
file_path = 'Data_Cuaca.xlsx'
data = pd.read_excel(file_path)

# 1. Mengambil kolom yang dibutuhkan
selected_columns = ['name', 'datetime', 'tempmax',
                    'tempmin']
data_selected = data[selected_columns]

# 2. Mengubah titik menjadi koma pada kolom tempmax dan
tempmin
data_selected['tempmax'] =
data_selected['tempmax'].str.replace(',',
    '.').astype(float)
data_selected['tempmin'] =
data_selected['tempmin'].str.replace(',',
    '.').astype(float)

# 3. Menghilangkan data duplikat
data_selected.drop_duplicates(inplace=True)

# 4. Menghilangkan baris dengan data kosong
data_selected.dropna(inplace=True)

# 5. Menghapus data outlier
data_selected = data_selected[(data_selected['tempmax']
    <= 40) & (data_selected['tempmin'] >= 20)]

# Menyimpan data yang sudah diolah ke file baru
(Optional)
output_file_path = 'Data_Cuaca_Bersih.xlsx'
data_selected.to_excel(output_file_path, index=False)

# Menampilkan data yang sudah diolah
print("Data Berhasil Dibersihkan")
```

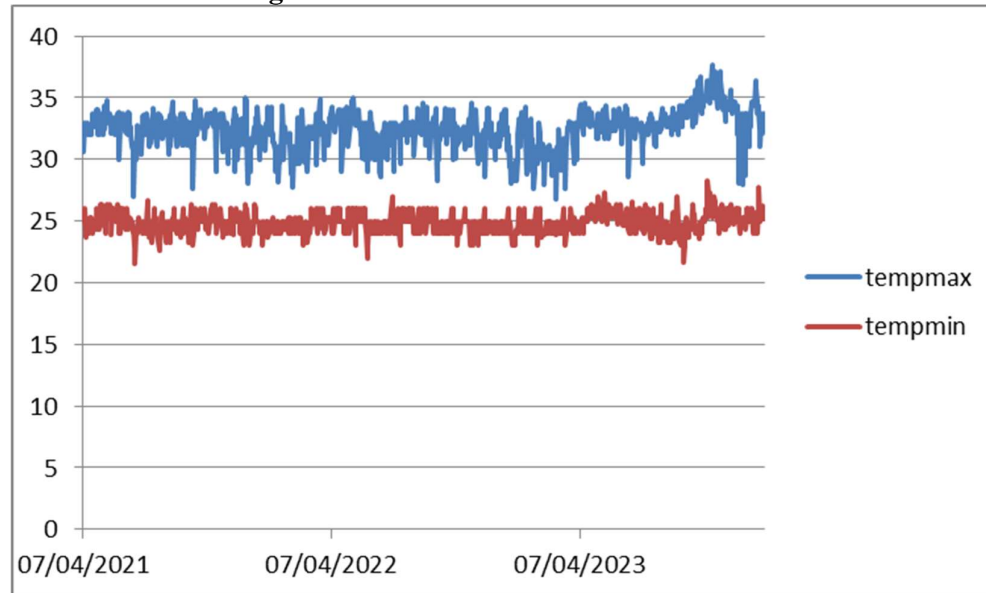
## 8.3 Output :

name	datetime	tempmax	tempmin
Jakarta	2021-04-07 00:00:00	30,6	25
Jakarta	2021-04-08 00:00:00	33	25
Jakarta	2021-04-09 00:00:00	32	26
Jakarta	2021-04-10 00:00:00	32	23,7
Jakarta	2021-04-11 00:00:00	32,7	24,6
Jakarta	2021-04-12 00:00:00	33	24,3
Jakarta	2021-04-13 00:00:00	33	25
Jakarta	2021-04-14 00:00:00	32	24,3
Jakarta	2021-04-15 00:00:00	32	24
Jakarta	2021-04-16 00:00:00	32	24,3
Jakarta	2021-04-17 00:00:00	32	25,3
Jakarta	2021-04-18 00:00:00	32	25,3
Jakarta	2021-04-19 00:00:00	32,7	25
Jakarta	2021-04-20 00:00:00	32,7	24
Jakarta	2021-04-21 00:00:00	33,7	24,7
Jakarta	2021-04-22 00:00:00	33	24,5
Jakarta	2021-04-23 00:00:00	32,7	24,3
Jakarta	2021-04-24 00:00:00	32,7	24,3



## BAB V VISUALISASI DATA

### 5.1. Time Series Pergerakan Suhu Maksimal dan Suhu Minimal



## BAB VI PEMROSESAN DATA MENGUNAKAN PYTHON

### 6.1. Pemrosesan Metode K-Means : sekaligus uji performa Cluster menggunakan metode Silhouette

#### 6.1.1. Modul Yang Digunakan :

- pandas
- matplotlib
- sklearn

#### 6.1.2. Kode Program :

```
import pandas as pd
from sklearn.cluster import KMeans
from sklearn.preprocessing import MinMaxScaler
from sklearn.metrics import silhouette_score
import matplotlib.pyplot as plt

# Membaca file Excel
file_path = "Data_Cuaca_Bersih.xlsx"
df = pd.read_excel(file_path)

# Memilih kolom tempmax dan tempmin
data = df[['tempmax', 'tempmin']]

# Scaling menggunakan Min-Max Scaling
scaler = MinMaxScaler()
data_scaled = scaler.fit_transform(data)

# Loop untuk mencoba nilai num_clusters dari 2 hingga 5
for num_clusters in range(2, 6):
    # Melakukan KMeans clustering
    kmeans = KMeans(n_clusters=num_clusters,
random_state=42)
    df['cluster'] = kmeans.fit_predict(data_scaled)

    # Menilai performa clustering menggunakan Silhouette
Score
    silhouette_avg = silhouette_score(data_scaled,
df['cluster'])
    print(f"Silhouette Score
(num_clusters={num_clusters}): {silhouette_avg}")

# Menyimpan hasil clustering ke dalam file Excel
output_file_path =
f"Data_Cuaca_Bersih_Clustered_{num_clusters}.xlsx"
df.to_excel(output_file_path, index=False)
```

```

    # Menampilkan dan menyimpan scatter plot hasil
    clustering
    plt.scatter(df['tempmax'], df['tempmin'],
c=df['cluster'], cmap='viridis')
    plt.title(f'Scatter Plot Hasil Clustering
(num_clusters={num_clusters})')
    plt.xlabel('TempMax')
    plt.ylabel('TempMin')
    plt.savefig(f'Scatter_Plot_Clustered_{num_clusters}.pn
g')
    plt.show()

```

## 6.2. Pemrosesan Metode Gaussian Mixture Models : sekaligus uji performa Cluster menggunakan metode Silhouette

### 6.2.1. Modul Yang Digunakan :

- pandas
- matplotlib
- statsmodels

### 6.2.2. Kode Program :

```

import pandas as pd
from sklearn.mixture import GaussianMixture
from sklearn.preprocessing import MinMaxScaler
from sklearn.metrics import silhouette_score
import matplotlib.pyplot as plt

# Membaca file Excel
file_path = "Data_Cuaca_Bersih.xlsx"
df = pd.read_excel(file_path)

# Memilih kolom tempmax dan tempmin
data = df[['tempmax', 'tempmin']]

# Scaling menggunakan Min-Max Scaling
scaler = MinMaxScaler()
data_scaled = scaler.fit_transform(data)

# Loop untuk mencoba nilai num_components dari 2 hingga 5
for num_components in range(2, 6):
    # Melakukan Gaussian Mixture Model clustering
    gmm = GaussianMixture(n_components=num_components,
random_state=42)
    df['cluster'] = gmm.fit_predict(data_scaled)

    # Menilai performa clustering menggunakan Silhouette
    Score

```

```
    silhouette_avg = silhouette_score(data_scaled,
df['cluster'])
    print(f"Silhouette Score
(num_components={num_components}): {silhouette_avg}")

    # Menyimpan hasil clustering ke dalam file Excel
    output_file_path =
f"Data_Cuaca_Bersih_Clustered_GMM_{num_components}.xlsx"
    df.to_excel(output_file_path, index=False)

    # Menampilkan dan menyimpan scatter plot hasil
clustering
    plt.scatter(df['tempmax'], df['tempmin'],
c=df['cluster'], cmap='viridis')
    plt.title(f'Scatter Plot Hasil Clustering GMM
(num_components={num_components})')
    plt.xlabel('TempMax')
    plt.ylabel('TempMin')
    plt.savefig(f'Scatter_Plot_Clustered_GMM_{num_componen
ts}.png')
    plt.show()
```

## BAB VII HASIL

### 7.1. Hasil Metode K-Means : sekaligus uji performa Cluster menggunakan metode Silhouette

#### 7.2.1. Output :

- Silhouette Score (num\_clusters=2): 0.37363226954587125
- Silhouette Score (num\_clusters=3): 0.3760170075621142
- Silhouette Score (num\_clusters=4): 0.35835217804018804
- Silhouette Score (num\_clusters=5): 0.3896178137614342

Berdasarkan hasil Silhouette Score yang telah diberikan, kita dapat menginterpretasikan kondisi rentang suhu cuaca sebagai berikut:

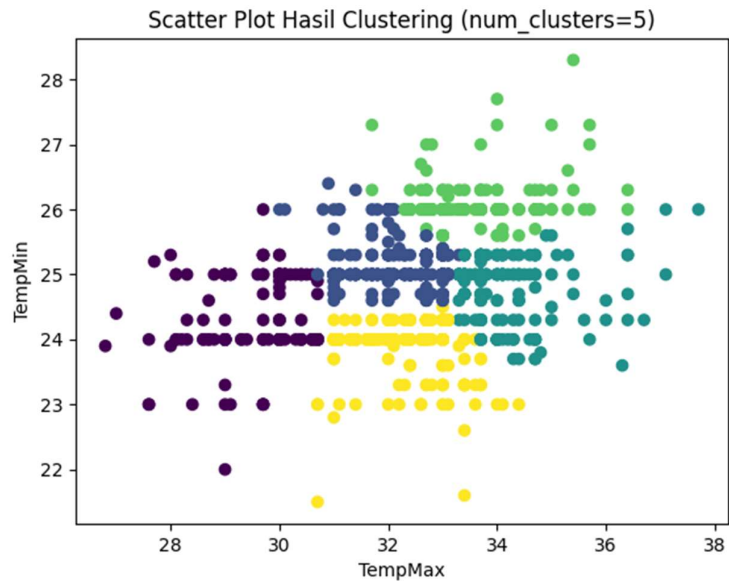
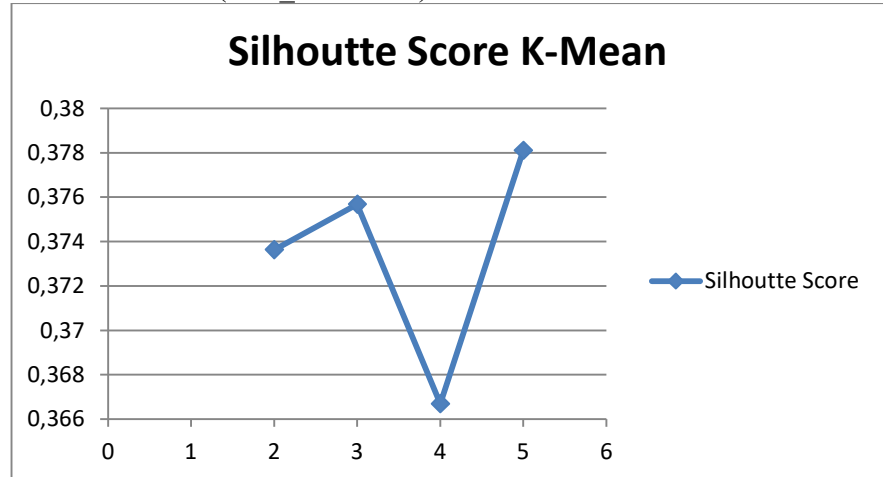
- **Variasi Ekstrem (num\_clusters=5):**
  - Silhouette Score: 0.3896
  - Pada konfigurasi dengan lima kluster, terdapat kecenderungan variasi ekstrem dalam rentang suhu cuaca. Nilai Silhouette Score yang relatif tinggi (0.3896) menandakan bahwa pembentukan lima kluster memberikan hasil yang cukup baik, dan mungkin terdapat variasi suhu yang signifikan antar kluster.
- **Kecenderungan Cuaca yang Konstan (num\_clusters=2, num\_clusters=3):**
  - Silhouette Score (num\_clusters=2): 0.3736
  - Silhouette Score (num\_clusters=3): 0.3760
  - Pada konfigurasi dengan dua dan tiga kluster, terdapat kemungkinan adanya kecenderungan cuaca yang konstan. Meskipun nilai Silhouette Score sedikit lebih rendah dibandingkan dengan lima kluster, namun tetap menunjukkan adanya struktur yang baik dalam pembentukan kluster.
- **Analisis Tambahan (num\_clusters=4):**
  - Silhouette Score: 0.3584
  - Konfigurasi dengan empat kluster memiliki nilai Silhouette Score yang lebih rendah, menunjukkan kemungkinan adanya beberapa kelompok suhu cuaca yang mungkin kurang terdefinisi dengan baik atau memiliki tumpang tindih.

Berdasarkan hasil tersebut, konfigurasi dengan lima kluster mungkin lebih menggambarkan variasi ekstrem dalam rentang

suhu cuaca, sementara konfigurasi dengan dua atau tiga kluster mengindikasikan adanya kecenderungan cuaca yang konstan.

### 7.2.2. Grafik :

- Silhouette Score (num\_clusters=5): 0.3896178137614342



Nilai Silhouette Score berkisar antara -1 hingga 1, dan semakin tinggi nilainya, semakin baik hasil clusteringnya. Dari hasil Uji performa Cluster K-Mean, bisa di lihat bahwa num\_clusters=5 yang mendapatkan skor Silhouette paling besar atau mendekati 1, menunjukkan bahwa pembentukan lima kluster memberikan hasil yang cukup baik. Nilai positif mengindikasikan bahwa objek-objek dalam kluster tersebut lebih serupa dengan anggota klusternya sendiri daripada dengan kluster lain.

## 7.2. Hasil Metode Gaussian Mixture Models : sekaligus uji performa Cluster menggunakan metode Silhouette

### 7.2.3. Output :

- SilhouetteScore (num\_components=2): 0.32422138482991714
- SilhouetteScore (num\_components=3): 0.32460686677351647
- Silhouett Score (num\_components=4): 0.23162912616883433
- Silhouete Score (num\_components=5): 0.16472131628154682

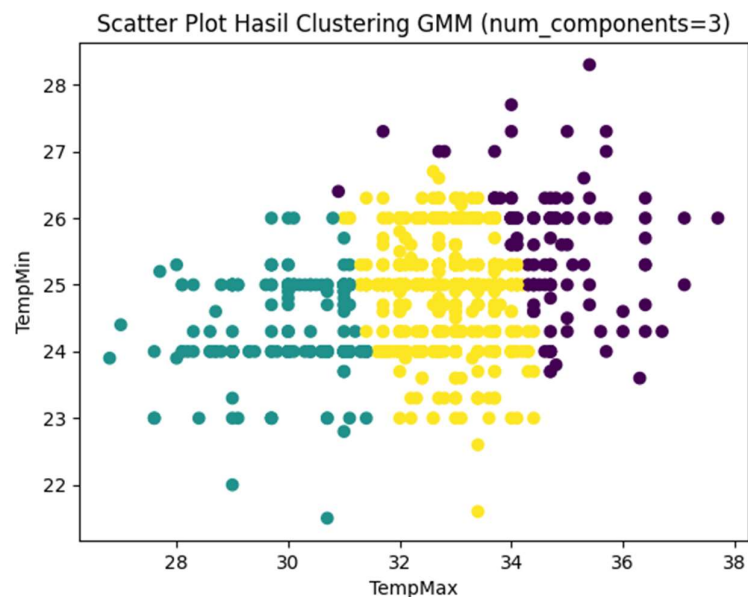
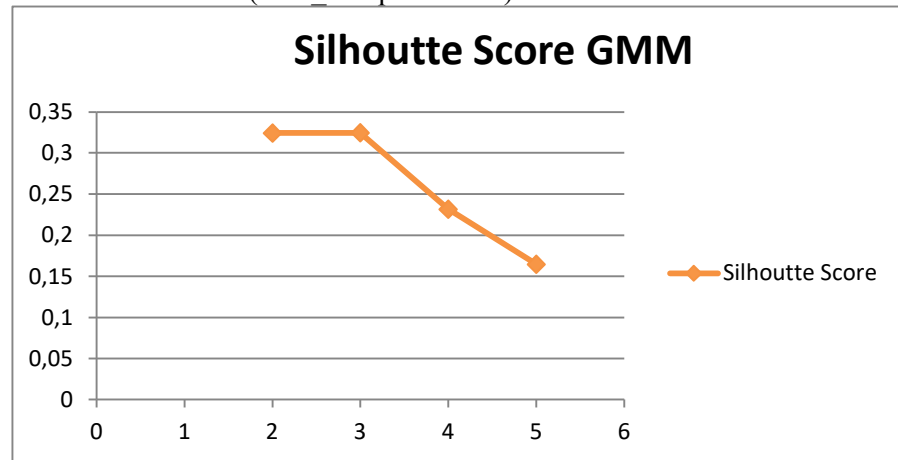
Berdasarkan hasil Silhouette Score, kita dapat menginterpretasikan kondisi rentang suhu cuaca sebagai berikut:

- **Variasi Ekstrem (num\_components=3 dan num\_components=2):**
  - Silhouette Score (num\_components=3): 0.3246
  - Silhouette Score (num\_components=2): 0.3242
  - Meskipun nilai Silhouette Score tidak sangat tinggi, konfigurasi dengan dua atau tiga komponen mungkin menunjukkan variasi ekstrem dalam rentang suhu cuaca. Keberadaan dua atau tiga komponen dapat mengindikasikan kecenderungan variasi suhu yang signifikan antar komponen.
- **Kecenderungan Cuaca yang Konstan (num\_components=4 dan num\_components=5):**
  - Silhouette Score (num\_components=4): 0.2316
  - Silhouette Score (num\_components=5): 0.1647
  - Konfigurasi dengan empat dan lima komponen memiliki nilai Silhouette Score yang lebih rendah, menunjukkan kemungkinan adanya kecenderungan cuaca yang konstan. Nilai yang lebih rendah mungkin menandakan adanya struktur yang kurang terdefinisi dengan baik atau tumpang tindih antar komponen.

Berdasarkan hasil tersebut, konfigurasi dengan dua atau tiga komponen mungkin lebih menggambarkan variasi ekstrem dalam rentang suhu cuaca, sementara konfigurasi dengan empat atau lima komponen mengindikasikan adanya kecenderungan cuaca yang konstan.

## 7.2.4. Grafik :

- SilhouetteScore (num\_components=3): 0.32460686677351647

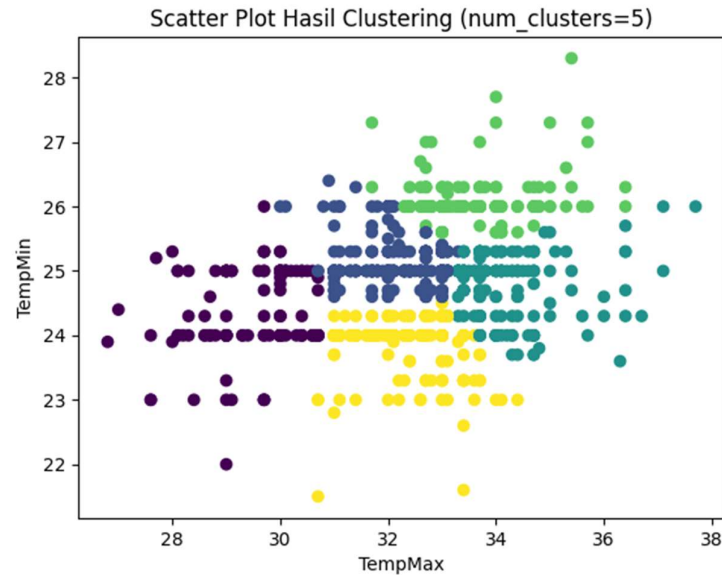


Dari hasil Uji performa Cluster Gaussian Mixture Models, bisa di lihat bahwa num\_component=3 yang mendapatkan skor Silhouette paling besar atau mendekati 1, menunjukkan bahwa pembentukan 3 kluster memberikan hasil yang cukup baik. Nilai positif mengindikasikan bahwa objek-objek dalam kluster tersebut lebih serupa dengan anggota klusternya sendiri daripada dengan kluster lain.

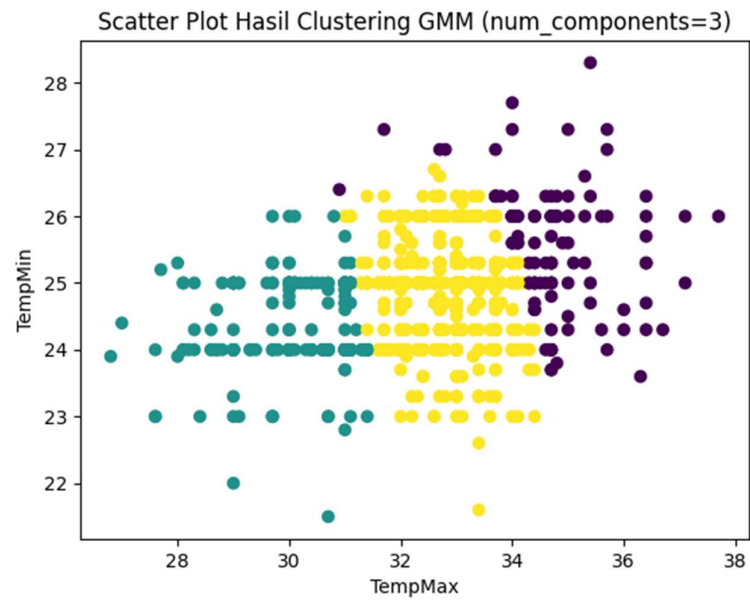


### 7.3. Perbandingan pemetaan data hasil metode K-Means dengan Gaussian Mixture Models

#### 7.3.1. Grafik K-Means:



#### 7.3.2. Grafik Gaussian Mixture Models:



Berdasarkan hasil Silhouette Score dan informasi kondisi rentang suhu cuaca, metode K-Means dengan 5 kluster dapat dianggap lebih optimal.

- **Lebih Tinggi Silhouette Score:** K-Means memiliki Silhouette Score yang lebih tinggi (0.3896) dibandingkan dengan Gaussian Mixture Models (GMM) (0.3246). Nilai Silhouette Score yang lebih tinggi menunjukkan kualitas klustering yang lebih baik.
- **Interpretasi Silhouette Score yang Lebih Tinggi:** Silhouette Score yang lebih tinggi pada metode K-Means menandakan bahwa objek-objek dalam kluster K-Means lebih serupa dengan anggota klusternya sendiri dan lebih berbeda dengan objek-objek dari kluster lain, memperlihatkan pembagian kelompok yang lebih baik.

## BAB VIII PENUTUP

### 8.1 Kesimpulan

8.1.1 Untuk membagi membagi data suhu cuaca harian berdasarkan nilai suhu minimal dan maksimal ke dalam kelompok-kelompok yang serupa dengan cara membuat kluster rentang kondisi suhu cuaca di jakarta berdasarkan nilai suhu cuaca minimal dan suhu cuaca maksimal jakarta tahun 2021-2023. berikut Langkah-langkahnya :

- a. Pemahaman data nilai suhu cuaca minimal dan suhu cuaca maksimal jakarta tahun 2021-2023
- b. Preprocessing terhadap data yang sudah disebutkan tadi
- c. Melakukan visualisasi Time Series Pergerakan Suhu Maksimal dan Suhu Minimal
- d. Kemudian, menerapkan algoritma K-Means dan GMM sekaligus uji performa Cluster menggunakan metode Silhouette per untuk mengelompokkan rentang kondisi suhu cuaca
- e. Menginterpretasikan kondisi cuaca per cluster berdasarkan hasil cluster suhu cuaca minimal dan suhu cuaca maksimal

Berikut insterpretasi hasil rentang kondisi suhu cuaca di jakarta berdasarkan nilai suhu cuaca minimal dan suhu cuaca maksimal jakarta tahun 2021-2023:

#### a. Metode K-Means

- **Variasi Ekstrem (num\_clusters=5):**

Pada konfigurasi dengan lima kluster, terdapat kecenderungan variasi ekstrem dalam rentang suhu cuaca. Nilai Silhouette Score yang relatif tinggi (0.3896) menandakan bahwa pembentukan lima kluster memberikan hasil yang cukup baik, dan mungkin terdapat variasi suhu yang signifikan antar kluster.

- **Kecenderungan Cuaca yang Konstan (num\_clusters=2, num\_clusters=3):**

Pada konfigurasi dengan dua dan tiga kluster, terdapat kemungkinan adanya kecenderungan cuaca yang konstan. Meskipun nilai Silhouette Score sedikit lebih rendah dibandingkan dengan lima kluster, namun tetap menunjukkan adanya struktur yang baik dalam pembentukan kluster.

- **Analisis Tambahan (num\_clusters=4):**

Konfigurasi dengan empat kluster memiliki nilai Silhouette Score yang lebih rendah, menunjukkan kemungkinan adanya beberapa kelompok suhu cuaca yang mungkin kurang terdefinisi dengan baik atau memiliki tumpang tindih.

#### b. Metode Gaussian Mixture Models

- **Variasi Ekstrem (num\_components=3 dan num\_components=2):**

Meskipun nilai Silhouette Score tidak sangat tinggi, konfigurasi dengan dua atau tiga komponen mungkin menunjukkan variasi ekstrem dalam rentang suhu cuaca. Keberadaan dua atau tiga komponen dapat mengindikasikan kecenderungan variasi suhu yang signifikan antar komponen.

- **Kecenderungan Cuaca yang Konstan (num\_components=4 dan num\_components=5):**

Konfigurasi dengan empat dan lima komponen memiliki nilai Silhouette Score yang lebih rendah, menunjukkan kemungkinan adanya kecenderungan cuaca yang konstan. Nilai yang lebih rendah mungkin menandakan adanya struktur yang kurang terdefinisi dengan baik atau tumpang tindih antar komponen.

8.1.2 Berdasarkan hasil Silhouette Score dan informasi kondisi rentang suhu cuaca, metode K-Means dengan 5 kluster dapat dianggap lebih optimal

- Lebih Tinggi Silhouette Score:** K-Means memiliki Silhouette Score yang lebih tinggi (0.3896) dibandingkan dengan Gaussian Mixture Models (GMM) (0.3246). Nilai Silhouette Score yang lebih tinggi menunjukkan kualitas klustering yang lebih baik.
- Interpretasi Silhouette Score yang Lebih Tinggi:** Silhouette Score yang lebih tinggi pada metode K-Means menandakan bahwa objek-objek dalam kluster K-Means lebih serupa dengan anggota klusternya sendiri dan lebih berbeda dengan objek-objek dari kluster lain, memperlihatkan pembagian kelompok yang lebih baik.