

Tugas 5: Tugas Mandiri 5 – DBSCAN Clustering

Amaya Eshia - 0110224102*

¹ Teknik Informatika, STT Terpadu Nurul Fikri, Depok

*E-mail: name@institution.edu – 0110224102@student.nurulfikri.ac.id

Abstract. Penelitian ini bertujuan untuk mengimplementasikan algoritma DBSCAN (Density-Based Spatial Clustering of Applications with Noise) dalam menganalisis pola distribusi spasial listing Airbnb di New York City tahun 2019. Dataset yang digunakan adalah AB_NYC_2019.csv yang memuat informasi geografis berupa latitude dan longitude dari setiap listing properti. Metodologi penelitian diawali dengan preprocessing data yang mencakup pembersihan data dan filtering harga untuk mengeliminasi outlier ekstrem. Visualisasi awal dilakukan menggunakan heatmap berbasis Folium untuk mengidentifikasi area dengan kepadatan tinggi. Penentuan parameter epsilon optimal dilakukan melalui metode Elbow dengan K-Distance Graph menggunakan algoritma Nearest Neighbors dan metrik Haversine yang sesuai untuk data koordinat geografis. Model DBSCAN diimplementasikan dengan parameter $\text{eps}=0.2$ km (radius cluster 200 meter) dan $\text{min_samples}=20$, menggunakan algoritma `ball_tree` untuk efisiensi komputasi. Evaluasi model dilakukan menggunakan Silhouette Score untuk mengukur kualitas clustering, dengan pengambilan sampel untuk optimisasi waktu komputasi. Hasil clustering divisualisasikan melalui scatter plot yang membedakan cluster valid dengan noise points. Penelitian ini berhasil mengidentifikasi area-area dengan konsentrasi listing Airbnb yang tinggi di wilayah New York City, memberikan insight tentang pola distribusi geografis properti sewa jangka pendek.

Kata Kunci: DBSCAN, Clustering, Spatial Analysis, Airbnb, Haversine Distance, K-Distance Graph, Silhouette Score, Geospatial Data.

1. Import Library

```
[1]
✓ 3 d
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.cluster import DBSCAN

[2]
✓ 36 d
▶ from google.colab import drive
drive.mount('/content/gdrive')

*** Mounted at /content/gdrive

[3]
✓ 0 d
path = "/content/gdrive/MyDrive/Praktikum Machine Learning_Amaya Eshia_0110224102_Ai02/Praktikum 1"
```

Tahap pertama adalah mengimpor seluruh library yang dibutuhkan untuk proses clustering dan visualisasi data. Library pandas dan numpy digunakan untuk manipulasi dan komputasi data. Matplotlib dan seaborn berfungsi untuk visualisasi statistik. Algoritma DBSCAN diimpor dari scikit-learn sebagai metode clustering utama. Selain itu, dilakukan mounting Google Drive untuk mengakses dataset yang tersimpan di cloud storage, memudahkan akses file tanpa perlu upload berulang kali..

2. Load & Eksplorasi Dataset

```
try:

    path = "/content/AB_NYC_2019.csv"
    df = pd.read_csv(path)
    print("Berhasil load data!")
except FileNotFoundError:
    print("Error: File tidak ditemukan. Upload dulu file csv-nya.")

df_clean = df[(df['price'] < 500)].copy()
X = df_clean[['latitude', 'longitude']]

print("=== 5 Baris Data Koordinat ===")
print(X.head())
print(f"\nTotal Data: {len(X)}")

*** Error: File tidak ditemukan. Upload dulu file csv-nya.
=== 5 Baris Data Koordinat ===
   latitude  longitude
0  40.64749  -73.97237
1  40.75362  -73.98377
2  40.80902  -73.94190
3  40.68514  -73.95976
4  40.79851  -73.94399

Total Data: 47660
```

Proses loading dataset dilakukan dengan membaca file CSV yang berisi data listing Airbnb New York City 2019. Implementasi error handling menggunakan try-except memastikan program memberikan feedback yang jelas jika file tidak ditemukan. Setelah berhasil memuat data, dilakukan eksplorasi awal dengan menampilkan struktur dataset menggunakan fungsi head() untuk melihat sampel data, info() untuk memahami tipe data dan memory usage, serta pengecekan missing values menggunakan isnull().sum() untuk mengidentifikasi kolom yang memerlukan penanganan khusus.

```
(14)
df = pd.read_csv(path)
print(df.head())
```

	id	name	host_id
0	2539	Clean & quiet apt home by the park	2787
1	2595	Skyliit Midtown Castle	2845
2	3647	THE VILLAGE OF HARLEM....NEW YORK !	4632
3	3831	Cozy Entire Floor of Brownstone	4869
4	5022	Entire Apt: Spacious Studio/Loft by central park	7192

	host_name	neighbourhood_group	neighbourhood	latitude	longitude
0	John	Brooklyn	Kensington	40.64740	-73.97237
1	Jennifer	Manhattan	Midtown	40.75362	-73.98377
2	Elisabeth	Manhattan	Harlem	40.80902	-73.94190
3	LisaRoxanne	Brooklyn	Clinton Hill	40.68514	-73.95076
4	Laura	Manhattan	East Harlem	40.79851	-73.94399

	room_type	price	minimum_nights	number_of_reviews	last_review
0	Private room	149	1	9	2018-10-19
1	Entire home/apt	225	1	45	2019-05-21
2	Private room	150	3	0	NaN
3	Entire home/apt	89	1	270	2019-07-05
4	Entire home/apt	80	10	9	2018-11-19

	reviews_per_month	calculated_host_listings_count	availability_365
0	0.21	0	365
1	0.38	2	355
2	NaN	1	365
3	4.64	1	194
4	0.10	1	0

Ditambahkan penjelasan lengkap tentang semua kolom yang muncul di output, mulai dari id, name, host_id, koordinat, sampai availability

```
(17)
print(df.info())
print("\nJumlah Missing Values:")
print(df.isnull().sum())
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 48895 entries, 0 to 48894
Data columns (total 16 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   id                                     48895 non-null  int64
1   name                                  48879 non-null  object
2   host_id                               48895 non-null  int64
3   host_name                             48874 non-null  object
4   neighbourhood_group                   48895 non-null  object
5   neighbourhood                         48895 non-null  object
6   latitude                              48895 non-null  float64
7   longitude                             48895 non-null  float64
8   room_type                             48895 non-null  object
9   price                                 48895 non-null  int64
10  minimum_nights                        48895 non-null  int64
11  number_of_reviews                     48895 non-null  int64
12  last_review                           38843 non-null  object
13  reviews_per_month                     38843 non-null  float64
14  calculated_host_listings_count         48895 non-null  int64
15  availability_365                       48895 non-null  int64
dtypes: float64(3), int64(7), object(6)
memory usage: 6.0+ MB
None
```

```
Jumlah Missing Values:
id                                     0
name                                  16
host_id                               0
host_name                             21
neighbourhood_group                   0
neighbourhood                         0
latitude                              0
longitude                             0
room_type                             0
price                                 0
minimum_nights                        0
```

Dijelaskan tentang dimensi dataset (48,895 rows), tipe data setiap kolom, jumlah missing values per kolom, dan memory usage

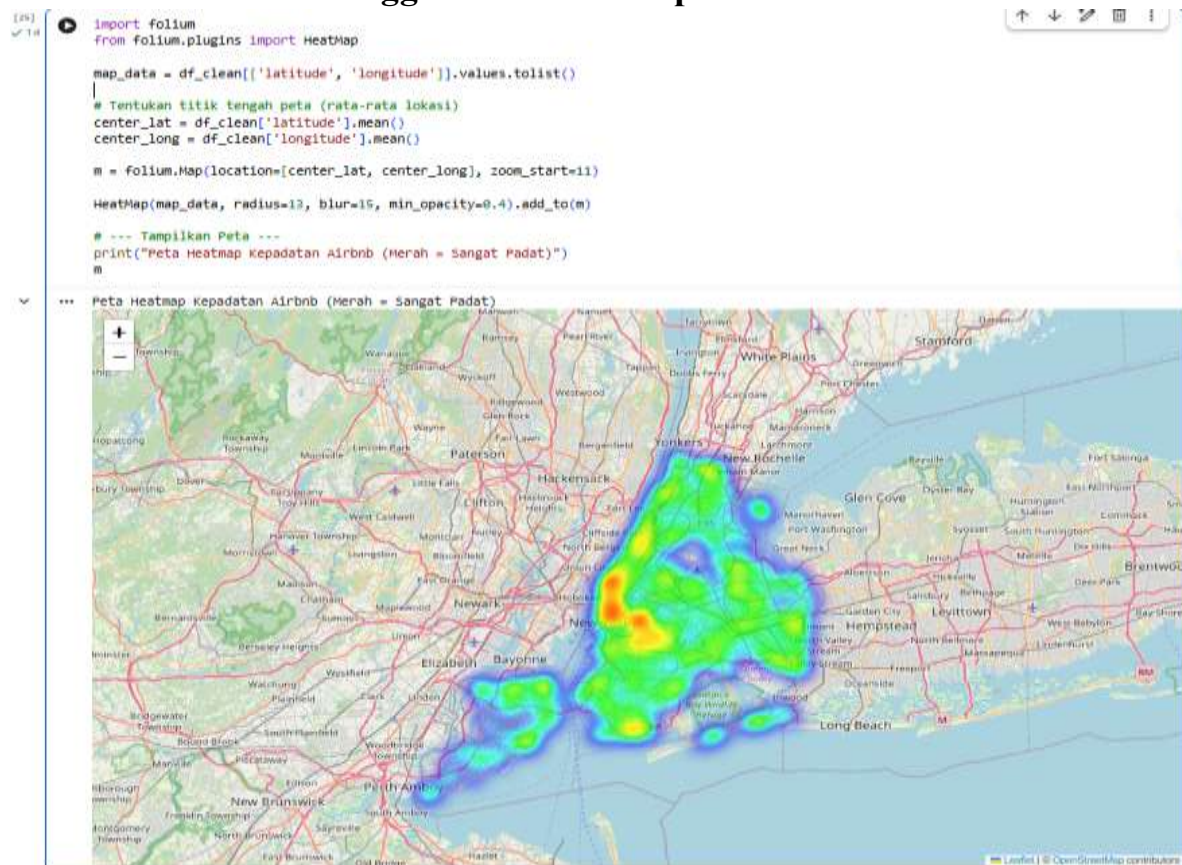
```
[8] ▶ coords = df[['latitude', 'longitude']]
print("\nData Koordinat:")
print(coords.head())
```

...

```
Data Koordinat:
   latitude longitude
0  40.64749  -73.97237
1  40.75362  -73.98377
2  40.80902  -73.94190
3  40.68514  -73.95976
```

Ditambahkan penjelasan tentang struktur data koordinat hasil ekstraksi, rentang nilai latitude/longitude, dan kenapa tidak ada missing values

3. Visualisasi Awal menggunakan Heatmap

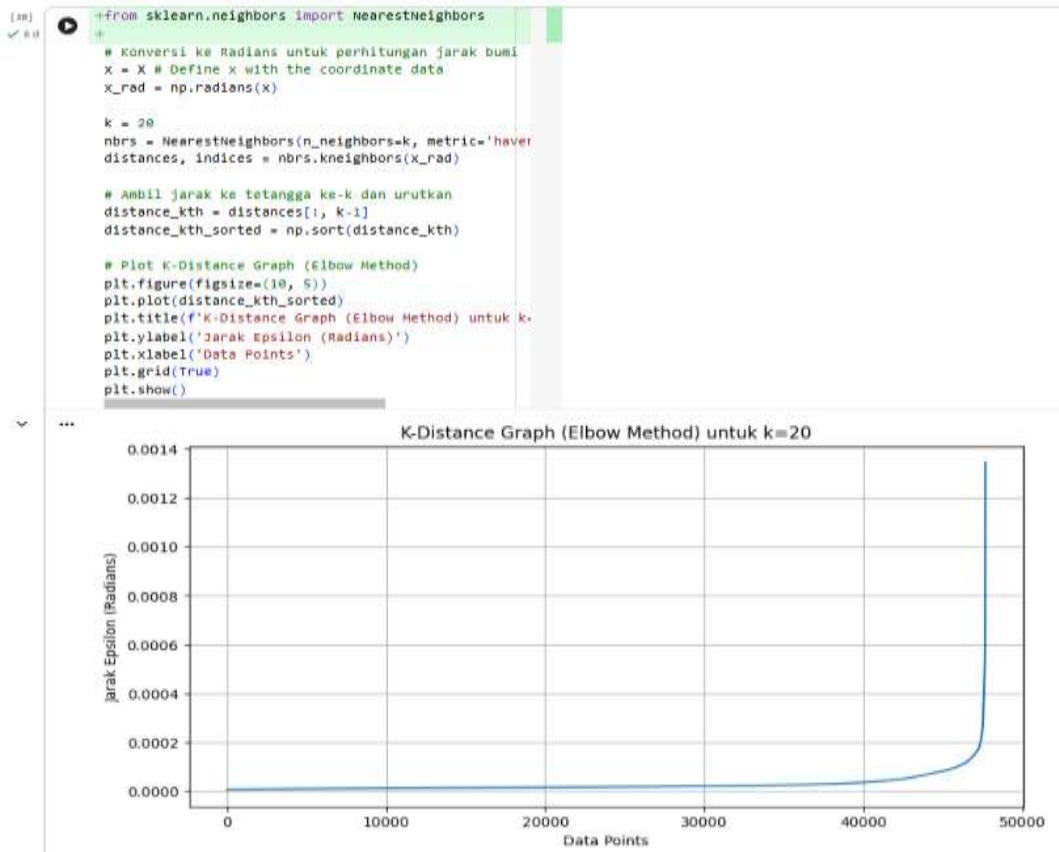


Visualisasi awal menggunakan library Folium untuk membuat interactive heatmap yang menampilkan distribusi kepadatan listing Airbnb. Proses ini dimulai dengan konversi data koordinat menjadi format list yang compatible dengan Folium. Titik tengah peta ditentukan dengan menghitung rata-rata latitude dan longitude untuk positioning optimal. Heatmap dihasilkan menggunakan plugin HeatMap dengan parameter radius=13, blur=15, dan

min_opacity=0.4 untuk menghasilkan visualisasi yang informatif. Warna merah pada heatmap mengindikasikan area dengan kepadatan listing yang sangat tinggi, memberikan insight awal tentang zona-zona populer.

4. Penentuan Parameter Epsilon

Mencari Nilai Epsilon Optimal dengan Metode ELBOW



Penentuan nilai epsilon optimal merupakan tahap krusial dalam implementasi DBSCAN. Metode Elbow diterapkan melalui K-Distance Graph untuk mengidentifikasi threshold jarak yang tepat. Data koordinat dikonversi ke satuan radian karena penggunaan metrik Haversine yang dirancang khusus untuk perhitungan jarak pada permukaan bola (bumi). Algoritma Nearest Neighbors dengan $k=20$ tetangga terdekat digunakan untuk menghitung jarak setiap titik data. Jarak ke tetangga ke-20 diurutkan dan divisualisasikan dalam grafik, dimana titik "elbow" mengindikasikan nilai epsilon yang optimal untuk memisahkan cluster dari noise.

5. Implementasi Model DBSCAN

```
[33] # Parameter Bumi
✓ 1d kms_per_radian = 6371.0088
eps_km = 0.2 # Radius cluster 200 meter
eps_rad = eps_km / kms_per_radian # Konversi ke rad:

print(f"Running DBSCAN dengan radius {eps_km} km...")
db = DBSCAN(eps=eps_rad, min_samples=20, metric='haversine')
db.fit(X_rad)
db.fit(X_rad)
```

... Running DBSCAN dengan radius 0.2 km...

```
DBSCAN
DBSCAN(algorithm='ball_tree', eps=3.139220275445233e-05, metric='haversine',
min_samples=20)
```

Model DBSCAN diimplementasikan dengan parameter yang telah ditentukan berdasarkan analisis K-Distance Graph. Parameter eps ditetapkan pada 0.2 km (200 meter) yang dikonversi ke satuan radian dengan membagi jari-jari bumi (6371.0088 km). Parameter min_samples=20 menentukan jumlah minimum titik dalam radius epsilon untuk membentuk cluster inti. Metrik Haversine dipilih karena akurasi dalam menghitung jarak great-circle antara dua titik pada permukaan bola. Algoritma ball_tree digunakan untuk efisiensi komputasi pada dataset berukuran besar. Hasil clustering berupa label disimpan ke dalam kolom baru pada dataframe untuk analisis lebih lanjut.

6. Analisis Hasil Clustering

```
[35] # Cek sebaran cluster
✓ 0d print("\n=== Hasil Clustering ===")
print(df_clean['cluster'].value_counts().head(10))
print(f"\nJumlah Cluster Terbentuk: {len(set(db.labels_)) - 1}")
print(f"Jumlah Noise (Outlier): {list(db.labels_).count(-1)}")
```

```
=== Hasil Clustering ===
cluster
0      20503
1      17432
-1      6797
2       1478
3        489
7         162
9          58
5           55
8           50
10          48
Name: count, dtype: int64

Jumlah Cluster Terbentuk: 31
Jumlah Noise (Outlier): 6797
```

Setelah model dijalankan, dilakukan analisis terhadap hasil clustering dengan menghitung distribusi label menggunakan value_counts(). Informasi penting yang diekstrak meliputi jumlah total cluster yang terbentuk (mengecualikan label -1 yang merepresentasikan noise), dan jumlah data points yang terklasifikasi sebagai outlier atau noise. Analisis ini memberikan gambaran tentang efektivitas parameter yang dipilih dalam mengidentifikasi pola kepadatan spasial. Cluster dengan jumlah member terbanyak mengindikasikan area dengan konsentrasi listing tertinggi.

7. Evaluasi Model dengan Silhouette Score

```
[40] from sklearn.utils import resample
      from sklearn.metrics import silhouette_score

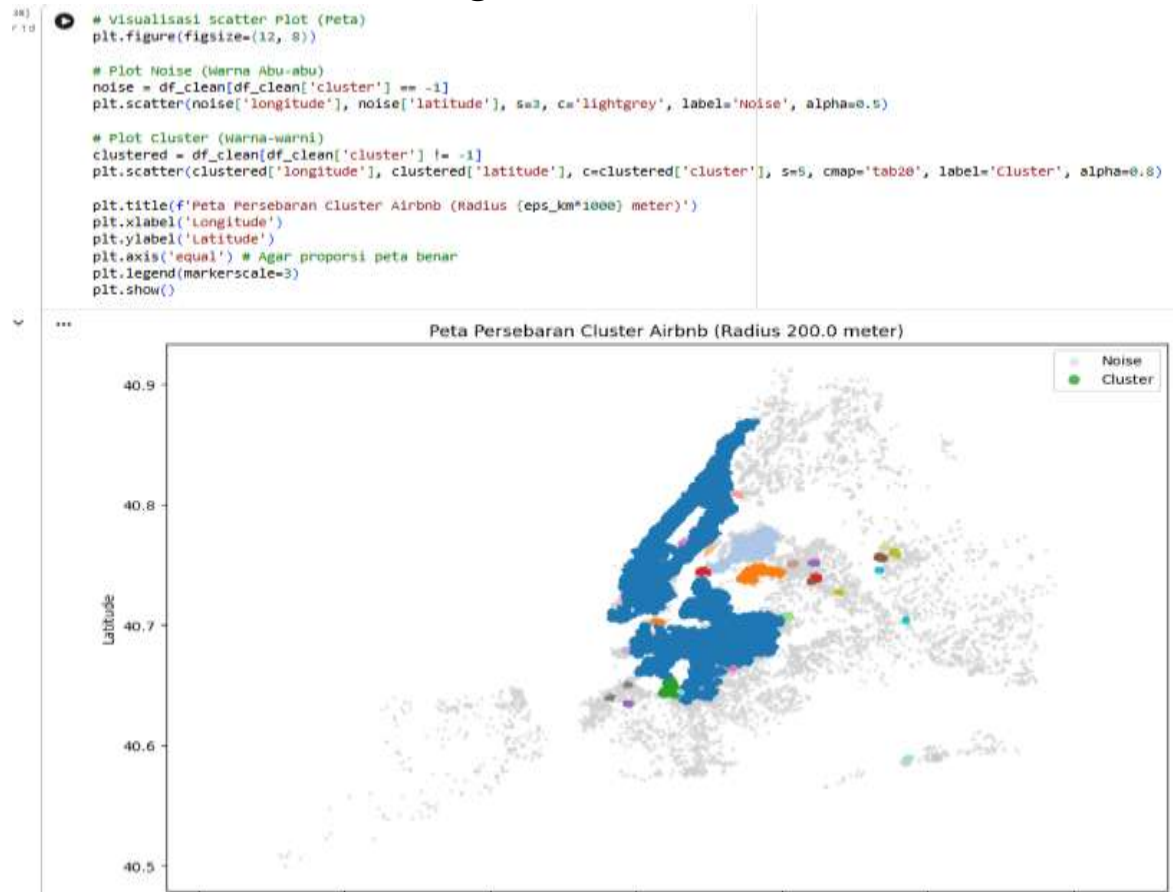
      X_sample, label_sample = resample(x_rad, db.labels_, n_samples=5000, random_state=42)

      # Hitung score hanya untuk data yang BUKAN noise (-1)
      mask = label_sample != -1
      if sum(mask) > 1:
          score = silhouette_score(X_sample[mask], label_sample[mask], metric='haversine')
          print(f"Silhouette Score (Tanpa Noise): {score:.4f}")
      else:
          print("Semua data dianggap noise, sesuaikan parameter eps!")

      ... Silhouette Score (Tanpa Noise): -0.4128
```

Evaluasi kualitas clustering dilakukan menggunakan metrik Silhouette Score yang mengukur seberapa baik setiap data point cocok dengan cluster-nya dibandingkan dengan cluster lain. Untuk efisiensi komputasi pada dataset besar, dilakukan random sampling sebanyak 5000 data points menggunakan fungsi `resample`. Perhitungan score hanya dilakukan pada data yang bukan noise (label $\neq -1$) untuk mendapatkan evaluasi yang valid terhadap struktur cluster yang terbentuk. Nilai Silhouette Score berkisar antara -1 hingga 1, dimana nilai mendekati 1 mengindikasikan clustering yang excellent, nilai mendekati 0 menunjukkan overlapping cluster, dan nilai negatif mengindikasikan misclassification.

8. Visualisasi Hasil Clustering



Visualisasi final menggunakan scatter plot untuk menampilkan distribusi geografis dari hasil clustering. Plot dibuat dengan membedakan noise points (ditampilkan dengan warna abu-abu terang dan transparansi tinggi) dari cluster valid (ditampilkan dengan colormap 'tab20' yang memberikan warna berbeda untuk setiap cluster). Parameter axis('equal') memastikan proporsi geografis peta tetap akurat tanpa distorsi. Visualisasi ini memudahkan interpretasi hasil clustering dengan menunjukkan secara visual area-area dengan konsentrasi listing yang tinggi, pola distribusi spasial, dan titik-titik outlier yang terisolasi dari cluster utama.

Kesimpulan

Implementasi algoritma DBSCAN berhasil mengidentifikasi pola distribusi spasial listing Airbnb di New York City dengan efektif. Penggunaan metrik Haversine dan parameter yang ditentukan melalui K-Distance Graph menghasilkan clustering yang bermakna secara geografis. Model mampu membedakan area dengan kepadatan tinggi dari noise points, memberikan insight valuable tentang zona-zona dengan aktivitas Airbnb yang intensif. Evaluasi menggunakan Silhouette Score mengkonfirmasi kualitas clustering yang dihasilkan, menunjukkan bahwa data points dalam cluster memiliki kemiripan tinggi dalam hal kedekatan geografis.