

Implementasi Algoritma K-means dan DBSCAN Tanpa Library Machine Learning

Nugroho Adi Susanto

Mei 2025

Contents

1 Pendahuluan

Dalam dokumen ini, kami mengimplementasikan algoritma K-means dan DBSCAN dari awal tanpa menggunakan library machine learning. Implementasi ini menggunakan dataset Iris untuk demonstrasi. Tujuan utama adalah memahami mekanisme internal dari kedua algoritma clustering populer ini serta membandingkan performa dan karakteristik keduanya.

2 Loading Dataset dan Persiapan Awal

Dataset Iris adalah salah satu dataset klasik dalam pembelajaran mesin. Dataset ini berisi 150 sampel dari tiga spesies Iris (Setosa, Versicolor, dan Virginica) dengan empat fitur: panjang sepal, lebar sepal, panjang petal, dan lebar petal.

2.1 Import Library yang Diperlukan

Berikut adalah library yang digunakan:

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn import datasets
6 from matplotlib.colors import ListedColormap
7 import random
8 from collections import defaultdict
9
10 # Set random seed untuk reproducibility
11 np.random.seed(42)
12 random.seed(42)
```

2.2 Loading Dataset Iris

Pada tahap ini, kita memuat dataset Iris dan menampilkan informasi dasar tentang dataset tersebut:

```
1 # Load dataset Iris
2 iris = datasets.load_iris()
3 X = iris.data # Fitur
4 y = iris.target # Label sebenarnya
5 feature_names = iris.feature_names
6 target_names = iris.target_names
7
8 # Membuat DataFrame untuk memudahkan visualisasi
9 iris_df = pd.DataFrame(X, columns=feature_names)
10 iris_df['species'] = [iris.target_names[i] for i in y]
11
12 # Tampilkan informasi dataset
13 print(f"Dataset shape: {X.shape}")
14 print(f"Feature names: {feature_names}")
15 print(f"Target names: {target_names}")
```

3 Visualisasi Data Awal

Sebelum menerapkan algoritma clustering, penting untuk memvisualisasikan data untuk mendapatkan pemahaman intuitif tentang distribusi data.

3.1 Visualisasi dengan Label Asli

Fungsi berikut digunakan untuk memvisualisasikan dataset Iris, dengan warna menunjukkan label asli atau hasil clustering:

```
1 # Fungsi untuk visualisasi dataset
2 def plot_iris(X, y, title):
3     plt.figure(figsize=(12, 5))
4
5     # Plot untuk 2 fitur pertama (sepal length vs sepal width)
6     plt.subplot(1, 2, 1)
7     plt.scatter(X[:, 0], X[:, 1], c=y, cmap=ListedColormap(['#FF0000', '#00FF00', '#0000FF']))
8     plt.xlabel(feature_names[0])
9     plt.ylabel(feature_names[1])
10    plt.title(f"{title} (Sepal Features)")
11
12    # Plot untuk 2 fitur terakhir (petal length vs petal width)
13    plt.subplot(1, 2, 2)
14    plt.scatter(X[:, 2], X[:, 3], c=y, cmap=ListedColormap(['#FF0000', '#00FF00', '#0000FF']))
15    plt.xlabel(feature_names[2])
16    plt.ylabel(feature_names[3])
17    plt.title(f"{title} (Petal Features)")
18
19    plt.tight_layout()
```

3.2 Visualisasi Data Tanpa Label

Dalam praktik clustering yang sebenarnya, kita biasanya tidak memiliki label yang benar dan bekerja dengan data tanpa label:

```
1 # Visualisasi data awal tanpa label (seperti yang akan kita lihat saat
   melakukan clustering)
2 plt.figure(figsize=(12, 5))
3
4 plt.subplot(1, 2, 1)
5 plt.scatter(X[:, 0], X[:, 1], c='gray')
6 plt.xlabel(feature_names[0])
7 plt.ylabel(feature_names[1])
8 plt.title("Iris Dataset tanpa Label (Sepal Features)")
9
10 plt.subplot(1, 2, 2)
11 plt.scatter(X[:, 2], X[:, 3], c='gray')
12 plt.xlabel(feature_names[2])
13 plt.ylabel(feature_names[3])
14 plt.title("Iris Dataset tanpa Label (Petal Features)")
```

4 Implementasi Algoritma K-means

K-means adalah algoritma clustering yang membagi data menjadi k kelompok berdasarkan jarak ke centroid terdekat. Berikut adalah implementasi K-means dari awal:

```
1 class KMeans:
2     def __init__(self, n_clusters=3, max_iterations=300, tolerance
3         =0.0001):
4         self.n_clusters = n_clusters
5         self.max_iterations = max_iterations
6         self.tolerance = tolerance
7         self.centroids = None
8         self.labels_ = None
9         self.inertia_ = None
10
11     def _initialize_centroids(self, X):
12         # Inisialisasi centroid secara acak dari data
13         idx = np.random.choice(X.shape[0], self.n_clusters, replace=
14             False)
15         return X[idx, :]
16
17     def _compute_distance(self, X, centroids):
18         # Menghitung jarak euclidean antara setiap titik dengan semua
19         centroid
20         distances = np.zeros((X.shape[0], self.n_clusters))
21         for k in range(self.n_clusters):
22             # Jarak euclidean kuadrat untuk mempercepat perhitungan
23             distances[:, k] = np.sum(np.square(X - centroids[k, :]),
24                 axis=1)
25         return distances
26
27     def _get_labels(self, distances):
28         # Menentukan cluster untuk setiap data berdasarkan jarak
29         terdekat
30         return np.argmin(distances, axis=1)
31
32     def _update_centroids(self, X, labels):
33         # Update centroid berdasarkan rata-rata data di setiap cluster
34         centroids = np.zeros((self.n_clusters, X.shape[1]))
35         for k in range(self.n_clusters):
36             if np.sum(labels == k) > 0: # Hindari pembagian dengan nol
37                 centroids[k, :] = np.mean(X[labels == k, :], axis=0)
38         return centroids
39
40     def _compute_inertia(self, X, labels, centroids):
41         # Menghitung inertia (jumlah kuadrat jarak data ke centroid
42         terdekat)
43         inertia = 0
44         for k in range(self.n_clusters):
45             if np.sum(labels == k) > 0:
46                 inertia += np.sum(np.square(X[labels == k, :] -
47                     centroids[k, :]))
48         return inertia
49
50     def fit(self, X):
51         # Inisialisasi centroid secara acak
52         self.centroids = self._initialize_centroids(X)
53         prev_centroids = np.copy(self.centroids)
```

```

47
48     for _ in range(self.max_iterations):
49         # Hitung jarak dan tentukan cluster
50         distances = self._compute_distance(X, self.centroids)
51         self.labels_ = self._get_labels(distances)
52
53         # Update centroid
54         self.centroids = self._update_centroids(X, self.labels_)
55
56         # Cek konvergensi
57         if np.all(np.abs(self.centroids - prev_centroids) < self.
tolerance):
58             break
59
60         prev_centroids = np.copy(self.centroids)
61
62         # Hitung inertia akhir
63         self.inertia_ = self._compute_inertia(X, self.labels_, self.
centroids)
64
65     return self
66
67     def predict(self, X):
68         # Prediksi cluster untuk data baru
69         distances = self._compute_distance(X, self.centroids)
70         return self._get_labels(distances)

```

4.1 Menjalankan K-means

Kita menjalankan K-means dengan k=3 (sesuai dengan jumlah spesies di dataset Iris):

```

1 # Implementasi K-means dengan k=3 (jumlah spesies di dataset Iris)
2 kmeans = KMeans(n_clusters=3)
3 kmeans.fit(X)
4
5 # Hasil clustering
6 y_kmeans = kmeans.labels_

```

4.2 Evaluasi Hasil K-means

Untuk mengevaluasi hasil clustering, kita menggunakan Adjusted Rand Index dan Silhouette Score:

```

1 from sklearn.metrics import adjusted_rand_score, silhouette_score
2
3 # Adjusted Rand Index - mengukur kesamaan antara dua pengelompokan
4 ari = adjusted_rand_score(y, y_kmeans)
5
6 # Silhouette Score - mengukur kualitas cluster yang terbentuk
7 silhouette = silhouette_score(X, y_kmeans)
8
9 print(f"Adjusted Rand Index: {ari:.4f}")
10 print(f"Silhouette Score: {silhouette:.4f}")

```

4.3 Elbow Method untuk K-means

Elbow Method digunakan untuk menentukan nilai k yang optimal. Metode ini melibatkan plot nilai inertia terhadap jumlah cluster:

```
1 # Implementasi Elbow Method untuk menentukan nilai k yang optimal
2 inertia_values = []
3 silhouette_values = []
4 k_values = range(2, 11)
5
6 for k in k_values:
7     kmeans = KMeans(n_clusters=k)
8     kmeans.fit(X)
9     inertia_values.append(kmeans.inertia_)
10    silhouette = silhouette_score(X, kmeans.labels_) if k > 1 else 0
11    silhouette_values.append(silhouette)
```

5 Implementasi Algoritma DBSCAN

DBSCAN (Density-Based Spatial Clustering of Applications with Noise) adalah algoritma clustering berbasis densitas yang dapat menemukan cluster dengan bentuk tidak beraturan dan mengidentifikasi noise points.

```
1 class DBSCAN:
2     def __init__(self, eps=0.5, min_samples=5):
3         self.eps = eps
4         self.min_samples = min_samples
5         self.labels_ = None
6
7     def _get_neighbors(self, X, point_idx):
8         # Mengembalikan indeks dari semua titik yang merupakan tetangga
9         # dari titik yang diberikan
10        distances = np.sum(np.square(X - X[point_idx, :]), axis=1)
11        return np.where(distances <= self.eps**2)[0]
12
13    def fit(self, X):
14        n_samples = X.shape[0]
15        # Inisialisasi label (-1 untuk noise, cluster diberi label 0, 1,
16        # 2, ...)
17        self.labels_ = np.full(n_samples, -1)
18
19        # Cluster saat ini
20        cluster = 0
21
22        # Untuk setiap titik di dataset
23        for point_idx in range(n_samples):
24            # Jika titik sudah diberi label, lanjut ke titik berikutnya
25            if self.labels_[point_idx] != -1:
26                continue
27
28            # Temukan tetangga dari titik saat ini
29            neighbors = self._get_neighbors(X, point_idx)
30
31            # Jika jumlah tetangga kurang dari min_samples, titik adalah
32            # noise
33            if len(neighbors) < self.min_samples:
```

```

31         self.labels_[point_idx] = -1 # Mark as noise
32         continue
33
34     # Titik adalah core point, mulai cluster baru
35     self.labels_[point_idx] = cluster
36
37     # Expand cluster: proses semua titik yang dapat dijangkau
    dari core point ini
38     # Inisialisasi seeds dengan tetangga dari core point
39     seeds = set(neighbors) - {point_idx}
40     seeds = list(seeds)
41
42     # Proses semua titik dalam seeds
43     seed_idx = 0
44     while seed_idx < len(seeds):
45         # Ambil titik dari seeds
46         current_point = seeds[seed_idx]
47
48         # Jika titik adalah noise, tambahkan ke cluster saat ini
49         if self.labels_[current_point] == -1:
50             self.labels_[current_point] = cluster
51
52         # Jika titik belum di-assign ke cluster apapun
53         if self.labels_[current_point] == -1:
54             # Tandai titik sebagai bagian dari cluster saat ini
55             self.labels_[current_point] = cluster
56
57         # Temukan tetangga dari titik saat ini
58         current_neighbors = self._get_neighbors(X,
    current_point)
59
60         # Jika titik adalah core point, tambahkan tetangga
    ke seeds
61         if len(current_neighbors) >= self.min_samples:
62             # Tambahkan tetangga yang belum diproses ke
    seeds
63             for neighbor in current_neighbors:
64                 if self.labels_[neighbor] == -1 or self.
    labels_[neighbor] == -1:
65                     if neighbor not in seeds:
66                         seeds.append(neighbor)
67
68             seed_idx += 1
69
70     # Lanjutkan dengan cluster selanjutnya
71     cluster += 1
72
73     return self

```

5.1 Menjalankan DBSCAN

Kita menjalankan DBSCAN dengan parameter awal:

```

1 # Menjalankan DBSCAN dengan parameter awal
2 dbscan = DBSCAN(eps=0.8, min_samples=5)
3 dbscan.fit(X)
4

```

```

5 # Hasil clustering dengan DBSCAN
6 y_dbscan = dbscan.labels_

```

5.2 Tuning Parameter DBSCAN

DBSCAN memerlukan dua parameter utama: `eps` (radius tetangga) dan `min_samples` (jumlah minimum tetangga). Kita perlu melakukan tuning untuk menemukan nilai optimal dari kedua parameter tersebut:

```

1 # Fungsi untuk menghitung jarak rata-rata ke k-tetangga terdekat
2 def compute_knn_distances(X, k):
3     # Menghitung jarak antara semua pasangan titik
4     distances = np.zeros((X.shape[0], X.shape[0]))
5     for i in range(X.shape[0]):
6         for j in range(X.shape[0]):
7             distances[i, j] = np.sqrt(np.sum(np.square(X[i] - X[j])))
8
9     # Untuk setiap titik, ambil jarak k-tetangga terdekat (kecuali
10    dirinya sendiri)
11    knn_distances = []
12    for i in range(X.shape[0]):
13        sorted_distances = np.sort(distances[i])
14        knn_distances.append(sorted_distances[k])
15
16    return np.array(knn_distances)
17
18 # Menghitung jarak rata-rata ke 5-tetangga terdekat
19 knn_distances = compute_knn_distances(X, 5)

```

5.3 Grid Search untuk Parameter DBSCAN

Grid search digunakan untuk secara sistematis mencoba berbagai kombinasi parameter dan memilih yang terbaik berdasarkan metrik evaluasi:

```

1 # Grid search untuk parameter DBSCAN
2 eps_values = np.linspace(0.2, 1.5, 14) # Range dari 0.2 hingga 1.5
3     dengan 14 nilai
4 min_samples_values = [3, 5, 7, 10, 15] # Nilai minPts yang akan dicoba
5
6 results = []
7
8 for eps in eps_values:
9     for min_samples in min_samples_values:
10        dbscan = DBSCAN(eps=eps, min_samples=min_samples)
11        dbscan.fit(X)
12
13        labels = dbscan.labels_
14        n_clusters = len(set(labels)) - (1 if -1 in labels else 0)
15        n_noise = list(labels).count(-1)
16
17        if n_clusters > 1: # Hanya hitung skor jika ada lebih dari satu
18            cluster
19            try:
20                silhouette = silhouette_score(X, labels)
21            except:
22                silhouette = float('nan')

```



```

21         ari = adjusted_rand_score(y, labels)
22     else:
23         silhouette = float('nan')
24         ari = float('nan')
25
26     results.append({
27         'eps': eps,
28         'min_samples': min_samples,
29         'n_clusters': n_clusters,
30         'n_noise': n_noise,
31         'silhouette': silhouette,
32         'ari': ari
33     })

```

5.4 Parameter Optimal DBSCAN

Setelah melakukan grid search, kita dapat menentukan parameter optimal untuk DBSCAN:

```

1  # Konversi hasil ke DataFrame
2  results_df = pd.DataFrame(results)
3
4  # Filter hasil dengan jumlah cluster > 1 dan tanpa nilai nan
5  valid_results = results_df[(results_df['n_clusters'] > 1) & (~results_df
6      ['silhouette'].isna())].copy()
7
8  if not valid_results.empty:
9      # Temukan parameter terbaik berdasarkan silhouette score
10     best_silhouette_idx = valid_results['silhouette'].idxmax()
11     best_params_silhouette = valid_results.loc[best_silhouette_idx]
12
13     # Temukan parameter terbaik berdasarkan ARI
14     best_ari_idx = valid_results['ari'].idxmax()
15     best_params_ari = valid_results.loc[best_ari_idx]

```

5.5 Menjalankan DBSCAN dengan Parameter Optimal

Kita menjalankan DBSCAN dengan parameter optimal yang ditemukan dari grid search:

```

1  # Menggunakan parameter terbaik dari grid search (berdasarkan silhouette
2  # score)
3  if 'best_params_silhouette' in locals():
4      optimal_eps = best_params_silhouette['eps']
5      optimal_min_samples = int(best_params_silhouette['min_samples'])
6  else: # Jika tidak ada hasil yang valid, gunakan default
7      optimal_eps = 0.6
8      optimal_min_samples = 5
9
10 # Jalankan DBSCAN dengan parameter optimal
11 optimal_dbscan = DBSCAN(eps=optimal_eps, min_samples=optimal_min_samples
12 )
13 optimal_dbscan.fit(X)
14 y_dbscan_optimal = optimal_dbscan.labels_

```

6 Perbandingan K-means dan DBSCAN

Pada bagian ini, kita membandingkan hasil clustering dari K-means dan DBSCAN:

```
1 # Visualisasi perbandingan hasil clustering K-means dan DBSCAN
2 plt.figure(figsize=(15, 10))
3
4 # Visualisasi label asli
5 plt.subplot(3, 2, 1)
6 plt.scatter(X[:, 0], X[:, 1], c=y, cmap=ListedColormap(['#FF0000', '#00
   FF00', '#0000FF']))
7 plt.xlabel(feature_names[0])
8 plt.ylabel(feature_names[1])
9 plt.title("Label Asli (Sepal Features)")
10
11 # Visualisasi hasil K-means
12 plt.subplot(3, 2, 3)
13 plt.scatter(X[:, 0], X[:, 1], c=y_kmeans, cmap=ListedColormap(['#FF0000',
   '#00FF00', '#0000FF']))
14 plt.xlabel(feature_names[0])
15 plt.ylabel(feature_names[1])
16 plt.title("Hasil K-means (Sepal Features)")
17
18 # Visualisasi hasil DBSCAN optimal
19 cmap_dbscan = ListedColormap(['#000000', '#FF0000', '#00FF00', '#0000FF',
   '#FFFF00', '#00FFFF', '#FF00FF'])
20 plt.subplot(3, 2, 5)
21 plt.scatter(X[:, 0], X[:, 1], c=y_dbscan_optimal, cmap=cmap_dbscan)
22 plt.xlabel(feature_names[0])
23 plt.ylabel(feature_names[1])
24 plt.title(f"Hasil DBSCAN (eps={optimal_eps:.2f}, minPts={
   optimal_min_samples}) (Sepal Features)")
```

7 Kesimpulan

Pada implementasi clustering dengan dataset Iris, kita telah melakukan perbandingan antara algoritma K-means dan DBSCAN dengan pengembangan dari awal (tanpa library machine learning).

7.1 Hasil K-means

- K-means berhasil dengan baik mengelompokkan dataset Iris menjadi 3 cluster, sesuai dengan jumlah spesies asli.
- Dengan Elbow Method, kita dapat mengonfirmasi bahwa $k=3$ memang merupakan nilai optimal untuk dataset ini.

7.2 Hasil DBSCAN

- DBSCAN memerlukan tuning parameter `eps` dan `min_samples` untuk mendapatkan hasil yang optimal.
- Setelah melakukan grid search, parameter terbaik ditemukan.

- DBSCAN memiliki keunggulan mampu mengidentifikasi noise points dan mendeteksi cluster dengan bentuk tidak beraturan.

7.3 Perbandingan

- K-means cenderung bekerja lebih baik pada dataset Iris karena cluster-nya bersifat konveks dan terpisah dengan baik.
- DBSCAN dapat menemukan cluster dengan bentuk yang lebih kompleks, tetapi memerlukan lebih banyak upaya dalam penentuan parameter.
- Evaluasi menggunakan Adjusted Rand Index dan Silhouette Score menunjukkan seberapa baik hasil clustering dibandingkan dengan label aslinya.

8 Referensi

1. Jain, A. K. (2010). Data clustering: 50 years beyond K-means. *Pattern Recognition Letters*, 31(8), 651-666.
2. Ester, M., Kriegel, H. P., Sander, J., & Xu, X. (1996, August). A density-based algorithm for discovering clusters in large spatial databases with noise. In *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining (KDD-96)* (Vol. 96, No. 34, pp. 226-231).
3. Scikit-learn: Machine Learning in Python, Pedregosa et al., *JMLR* 12, pp. 2825-2830, 2011.