

Machine Learning

Tugas 2



Agung Nursatria Banyuwiguna

1301150073

IF-39-03

Program Studi S1 Teknik Informatika – Fakultas Informatika

Universitas Telkom

Jalan Telekomunikasi Terusan Buah Batu, Bandung

Indonesia

1. Analisis Masalah

Diketahui terdapat 688 titik yang terpengar dalam grafik/plot 2 dimensi dengan tanpa label. Dibutuhkan suatu pengelompokkan titik - titik tersebut secara optimal dengan menggunakan metode k-means. Untuk mendapatkan pengelompokkan yang optimal, dibutuhkan centroid dengan jumlah tepat untuk menjadi titik tengah dari kelompok yang ada. Jumlah tepat centroid ditentukan berdasarkan SSE yang merupakan jarak dari anggota ke centroid. Digunakan metode elbow yang akan menyatakan jumlah centroid optimal dimana k turun drastis terakhir. Disini penulis menyatakan turun drastis jika penurunan SSE setidaknya lebih dari 30% penurunan pertama.

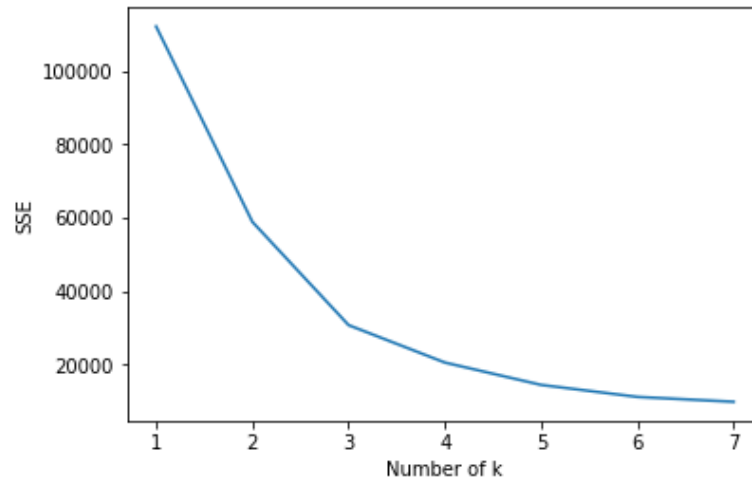
2. Desain System

Berikut langkah untuk membangun sistem k-means:

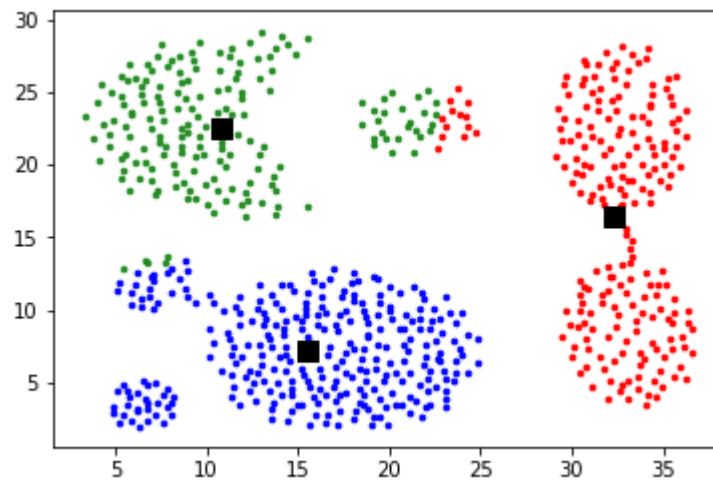
1. Load data dari file
2. Cari K optimal dengan method elbow
3. Inisialisasi centroid sejumlah K
4. Lakukan clustering menentukan kelompok dari titik-titik yang ada berdasarkan jarak ke centroid
5. Inisialisasi centroid baru
6. ulangi langkah 4 dan 5 hingga centroid lama dan centroid baru ditempat yang sama
7. Keluarkan hasil
8. Simpan hasil clustering [opsional]

3. Evaluasi Hasil Running

Dengan data training:



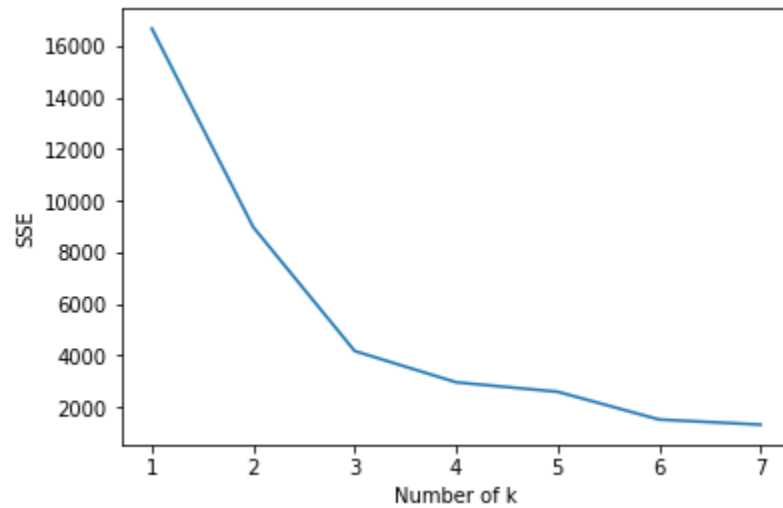
Berdasarkan method elbow pada grafik diatas, ditentukan k optimal yaitu 3. Dikarenakan setelah 3, penurunan tidak drastis lagi.



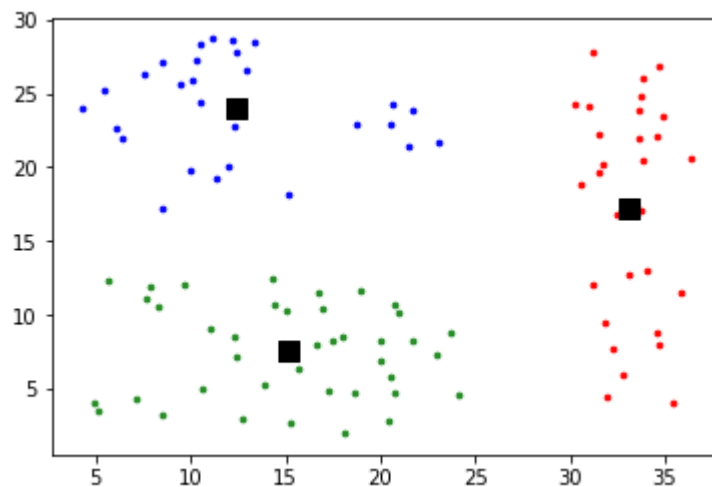
Dengan k=3, didapat hasil clustering seperti gambar diatas dengan perbandingan masing – masing titik:

```
cluster
total   : 688 titik
0       : 216 titik
1       : 294 titik
2       : 178 titik
```

Dengan data testing:



Berdasarkan method elbow pada grafik diatas, ditentukan k optimal yaitu 3. Dikarenakan setelah 3, penurunan tidak drastis lagi.



Dengan $k=3$, didapat hasil clustering seperti gambar diatas dengan perbandingan masing – masing titik:

```
cluster
total   : 100 titik
0       : 42 titik
1       : 28 titik
2       : 30 titik
```

4. Screenshot program

Main.py

```

17 import numpy as np
18 from matplotlib import pyplot as plt
19 from elbow import optimal_k
20 from kmeans import initCentroid, clustering
21
22 # Load data dan jadikan array
23 data = np.loadtxt('TrainsetTugas2.txt')
24 x_data = data[:,0]
25 y_data = data[:,1]
26
27 k= optimal_k(data)
28
29 # koordinat X,Y centroid data random
30 centroid = initCentroid(data,k)
31
32 # clustering
33 clusters = clustering(data,centroid,k)
34
35 # Buat Plot
36 colors = ['blue', 'forestgreen', 'red', 'darkgoldenrod', 'purple', 'cyan']
37 fig, ax = plt.subplots()
38 print('cluster')
39 print('total\t:', len(data), 'titik')
40 for i in range(k):
41     titik = np.array([data[j] for j in range(len(data)) if clusters[j] == i])
42     print(i, '\t:', len(titik), 'titik')
43     ax.scatter(titik[:, 0], titik[:, 1], s=7, c=colors[i])
44 ax.scatter(centroid[:, 0], centroid[:, 1], marker='s', s=100, c='black')
45
46 # simpan result ke txt
47 result = []
48 result = np.column_stack((data, clusters))
49 np.savetxt('result_trains.txt', result, newline='\n', fmt="%s\t%s\t%i")

```

Test.py

```

18 import numpy as np
19 from matplotlib import pyplot as plt
20 from elbow import optimal_k
21 from kmeans import initCentroid, clustering
22
23 # Load data dan jadikan array
24 data = np.loadtxt('TestsetTugas2.txt')
25 x_data = data[:,0]
26 y_data = data[:,1]
27
28 k= optimal_k(data)
29
30 # koordinat X,Y centroid data random
31 centroid = initCentroid(data,k)
32
33 # clustering
34 clusters = clustering(data,centroid,k)
35
36 # Buat Plot
37 colors = ['blue', 'forestgreen', 'red', 'darkgoldenrod', 'purple', 'cyan']
38 fig, ax = plt.subplots()
39 print('cluster')
40 print('total\t:', len(data), 'titik')
41 for i in range(k):
42     titik = np.array([data[j] for j in range(len(data)) if clusters[j] == i])
43     print(i, '\t:', len(titik), 'titik')
44     ax.scatter(titik[:, 0], titik[:, 1], s=7, c=colors[i])
45 ax.scatter(centroid[:, 0], centroid[:, 1], marker='s', s=100, c='black')
46
47 # simpan result ke txt
48 result = []
49 result = np.column_stack((data, clusters))
50 np.savetxt('result_test.txt', result, newline='\n', fmt="%s\t%s\t%i")

```

Euclidean.py

```
8 import numpy as np
9
10 # Euclidean distance untuk banyak data centroid
11 def euclid(data, centroid, ax=1):
12     return np.linalg.norm(data - centroid, axis=ax)
```

Kmeans.py

```
8 import numpy as np
9 from copy import deepcopy
10 from euclidean import euclid
11
12 def initCentroid(data,k):
13     x_centroid = np.random.randint(0, np.max(data), size=k)
14     y_centroid = np.random.randint(0, np.max(data), size=k)
15     centroid = np.array(list(zip(x_centroid, y_centroid)), dtype=np.float32)
16     return centroid
17
18 def clustering(data, centroid, k):
19     old_centroid = np.zeros(centroid.shape)
20     centers = np.zeros(len(data))
21     sameLocation = euclid(centroid, old_centroid, None)
22     while sameLocation != 0:
23         for i in range(len(data)):
24             distances = euclid(data[i], centroid)
25             center = np.argmin(distances)
26             centers[i] = center
27         old_centroid = deepcopy(centroid)
28         for i in range(k):
29             titik = [data[j] for j in range(len(data)) if centers[j] == i]
30             centroid[i] = np.mean(titik, axis=0)
31         sameLocation = euclid(centroid, old_centroid, None)
32     return centers
```

Elbow.py

```
8 from kmeans import clustering,initCentroid
9 from matplotlib import pyplot as plt
10
11 def optimal_k(data):
12     # mencari daftar sse
13     sse = []
14     for k in range(1,8):
15         centroid = initCentroid(data,k)
16         clusters = clustering(data,centroid,k)
17         sumtoCluster = 0
18         for i in range(k):
19             for j in range(len(data)):
20                 if clusters[j] == i:
21                     # Euclidean Distance untuk sumtoCluster
22                     sumtoCluster += (abs((data[j,0]-centroid[i][0])**2) +
23                                     abs((data[j,1]-centroid[i][1])**2))
24             sse.append([k,sumtoCluster])
25
26     # Plot grafik kluster
27     sseKey = [sse[j][0] for j in range(len(sse))]
28     sseValue = [sse[j][1] for j in range(len(sse))]
29     plt.figure()
30     plt.plot(sseKey, sseValue)
31     plt.xlabel("Number of k")
32     plt.ylabel("SSE")
33     plt.show()
34
35     # mencari k optimal dengan patokan jika nilai yang berkurang lebih samu
36     k_optimal = 2
37     first_drop_sse = (sse[0][1]-sse[1][1])
38     for k in range(1,5):
39         if (sse[k][1]-sse[k+1][1]) >= (first_drop_sse*0.3)):
40             k_optimal = sse[k+1][0]
41     return k_optimal
```