

K-Means

1. Cara kerja algoritma K-Means

K-Means adalah algoritma clustering yang bertujuan untuk membagi data menjadi sejumlah k cluster yang ditentukan sebelumnya, berdasarkan kedekatan jarak antara titik data dan centroid cluster. Algoritma ini bekerja dengan iteratif untuk meminimalkan variasi dalam cluster (intra-cluster variance). Langkah-langkah utama dalam algoritma K-Means adalah sebagai berikut.

Inisialisasi:

- Pilih jumlah cluster (k) yang diinginkan.
- Inisialisasi centroid awal. Ini bisa dilakukan secara acak atau dengan metode k-means++ yang memilih centroid lebih cerdas untuk memperbaiki konvergensi.

Assignment Step (Penugasan):

- Untuk setiap titik data, hitung jarak antara titik data tersebut dan setiap centroid.
- Tugaskan titik data ke cluster yang centroid-nya paling dekat (dengan jarak terkecil).

Update Step (Pembaruan):

- Setelah semua titik data ditugaskan ke cluster, hitung ulang posisi centroid dengan mengambil rata-rata dari semua titik data yang berada dalam cluster tersebut.

Iterasi:

- Ulangi langkah Assignment dan Update sampai centroid tidak berubah lagi (konvergensi), atau hingga jumlah iterasi maksimum tercapai.

Hasil Akhir:

- Algoritma berhenti, dan hasil akhir adalah pembagian data menjadi k cluster dengan centroid yang optimal berdasarkan minimisasi variasi dalam cluster.

4. Hasil Evaluasi Model

From Scratch	From Scikit-Learn
<pre>Centroids: [[4.95588982e-01 4.34131737e-02 7.93413174e-02 2.57485030e-01 4.35628743e-01 1.09266467e+01 4.68113772e+00 4.79041916e-01 6.79640719e-01 2.30089820e+00 1.13772455e-01 1.00299401e-01 2.12574850e-01 -1.13203694e-01 3.50337050e-02 -2.06582594e-01 2.31836751e-01] [4.82315113e-01 7.60986066e-02 5.68060021e-02 1.77920686e-01 4.94105038e-01 7.66988210e+00 4.55948553e+00 4.19078242e-01 7.68488746e-01 2.21757771e+00 1.56484459e-01 4.60878885e-02 9.53912111e-02 1.10365889e-01 1.02513925e-01 -3.23758597e-02 -5.50647916e-02] [3.42751843e-01 7.37100737e-02 1.10565111e-02 5.15970516e-02 5.19656020e-01 2.67444717e+00 4.15970516e+00 1.31449631e-01 8.37837838e-01 2.10196560e+00 1.62162162e-01 1.84275184e-02 9.82800983e-03 -3.36011145e-02 -1.46250623e-01 2.06638636e-01 -1.27139434e-01]] Silhouette Score: 0.2059778191835666 Davies-Bouldin Index: 1.6515535190819286</pre>	<pre>Centroids: [[4.96327387e-01 4.51206716e-02 7.45015740e-02 2.33997901e-01 4.43861490e-01 1.03504722e+01 4.66946485e+00 4.81636936e-01 6.97796432e-01 2.28961175e+00 1.08079748e-01 8.81427072e-02 1.85729276e-01 -6.62628641e-02 5.38594792e-02 -1.84498996e-01 1.91931141e-01] [4.45121951e-01 8.41463415e-02 5.00000000e-02 1.58536585e-01 5.04878049e-01 6.64756098e+00 4.43414634e+00 3.48780488e-01 7.96341463e-01 2.11707317e+00 1.78048780e-01 3.90243902e-02 7.31707317e-02 1.52477468e-01 8.71993049e-02 4.09587821e-02 -1.30459479e-01] [3.45794393e-01 7.47663551e-02 4.67289720e-03 4.20560748e-02 5.26479751e-01 2.05140187e+00 4.17601246e+00 1.13707165e-01 8.33333333e-01 2.17912773e+00 1.63551402e-01 1.40186916e-02 3.11526480e-03 -9.63909881e-02 -1.91326345e-01 2.21559722e-01 -1.18276643e-01]] Silhouette Score: 0.22276989899107055 Davies-Bouldin Index: 1.5664437376513092</pre>

Hasil centroid dari implementasi K-Means scratch menunjukkan distribusi centroid yang berbeda dari yang dihasilkan oleh sklearn. Hal ini mungkin disebabkan oleh perbedaan dalam metode inisialisasi centroid, cara pengukuran jarak, atau konvergensi yang digunakan dalam implementasi tersebut.

Model K-Means yang diimplementasikan from scratch menghasilkan Silhouette Score sebesar 0.206 dan Davies-Bouldin Index sebesar 1.652. Silhouette Score yang rendah menunjukkan bahwa jarak antara data dalam cluster relatif dekat, tetapi jarak antara data di cluster yang berbeda cukup jauh, namun masih rendah, menunjukkan bahwa hasil clustering tidak terlalu baik. Davies-Bouldin Index yang lebih tinggi mengindikasikan bahwa rata-rata rasio jarak intra-cluster terhadap inter-cluster adalah lebih tinggi, yang berarti ada ruang untuk perbaikan dalam segregasi cluster.

Sebaliknya, model K-Means yang diimplementasikan menggunakan scikit-learn menghasilkan Silhouette Score sebesar 0.223 dan Davies-Bouldin Index sebesar 1.566. Meskipun Silhouette Score-nya masih rendah, ia menunjukkan sedikit peningkatan dibandingkan dengan implementasi from scratch, dan Davies-Bouldin Index yang lebih rendah menunjukkan bahwa model scikit-learn sedikit lebih baik dalam mengelompokkan data dengan membedakan cluster.

Hal tersebut terjadi karena implementasi from scratch mungkin memiliki keterbatasan dalam hal optimasi, seperti konvergensi centroid atau penanganan kasus khusus, yang dapat mempengaruhi hasil akhir. Metode inisialisasi centroid juga dapat mempengaruhi performa model; dalam hal ini, scikit-learn menggunakan metode K-Means++ yang lebih canggih untuk inisialisasi centroid, yang secara umum memberikan hasil yang lebih stabil dan lebih baik dibandingkan dengan inisialisasi acak atau sederhana.

5. Improvement

Untuk mencapai hasil clustering yang lebih baik, penting untuk melakukan optimasi konvergensi algoritma K-Means dengan memastikan bahwa centroid benar-benar mencapai posisi stabil dan tidak bergerak secara signifikan antara iterasi. Penanganan kasus khusus juga dapat dilakukan dengan memeriksa dan menangani kasus khusus seperti cluster kosong atau fitur dengan variansi sangat kecil yang mungkin mempengaruhi hasil clustering.

Selain itu, evaluasi hasil clustering harus dilakukan secara menyeluruh dengan mempertimbangkan untuk menggunakan berbagai metode tambahan seperti Dunn Index dan Elbow Method. Jika memungkinkan, lakukan cross-validation pada data untuk memastikan bahwa hasil clustering konsisten dan tidak bergantung pada subset data tertentu. Analisis mendalam terhadap cluster yang dihasilkan juga dapat membantu memastikan bahwa hasil clustering mencerminkan struktur data yang diharapkan dan memberikan wawasan tambahan tentang efektivitas model.