# Unsupervised Learning: Concept & Applications

Unsupervised learning finds hidden structures in data without relying on ground-truth labels. Algorithms identify patterns — such as groupings, low-dimensional representations, or rare anomalies — using only the input features.

## Key Applications

### Clustering (Group Discovery)

Organise large image corpuses by visual similarity for efficient pre-labelling or exploration.

### Anomaly Detection (Rare-Event Discovery)

Identify potential network intrusions from connection features, where outliers indicate security threats.

# Clustering & K-Means

## K-Means: objective, algorithm, and properties

The K-Means algorithm is a fundamental method in unsupervised learning for partitioning data.

### Goal

Partition (n) samples into (k) clusters so each sample is assigned to the nearest cluster centroid.

### Objective (WCSS)

$$J = \sum_{1-i}^{k} \sum_{x \in C_i} |x - \mu_i|^2$$

Minimize (J) (tight intra-cluster variance).

### Algorithm (iterative coordinate descent)

Initialization: choose (k) centroids (random, K-Means++).

Update: $\mu_i = \frac{1}{|C_i|}\sum_{x \in C_i} x$

**1** ——— **2** ——— **3** ——— **4**

Assignment: $c(\mathbf{x}) = \arg\min_{j}|\mathbf{x} - \boldsymbol{\mu}_j|$.

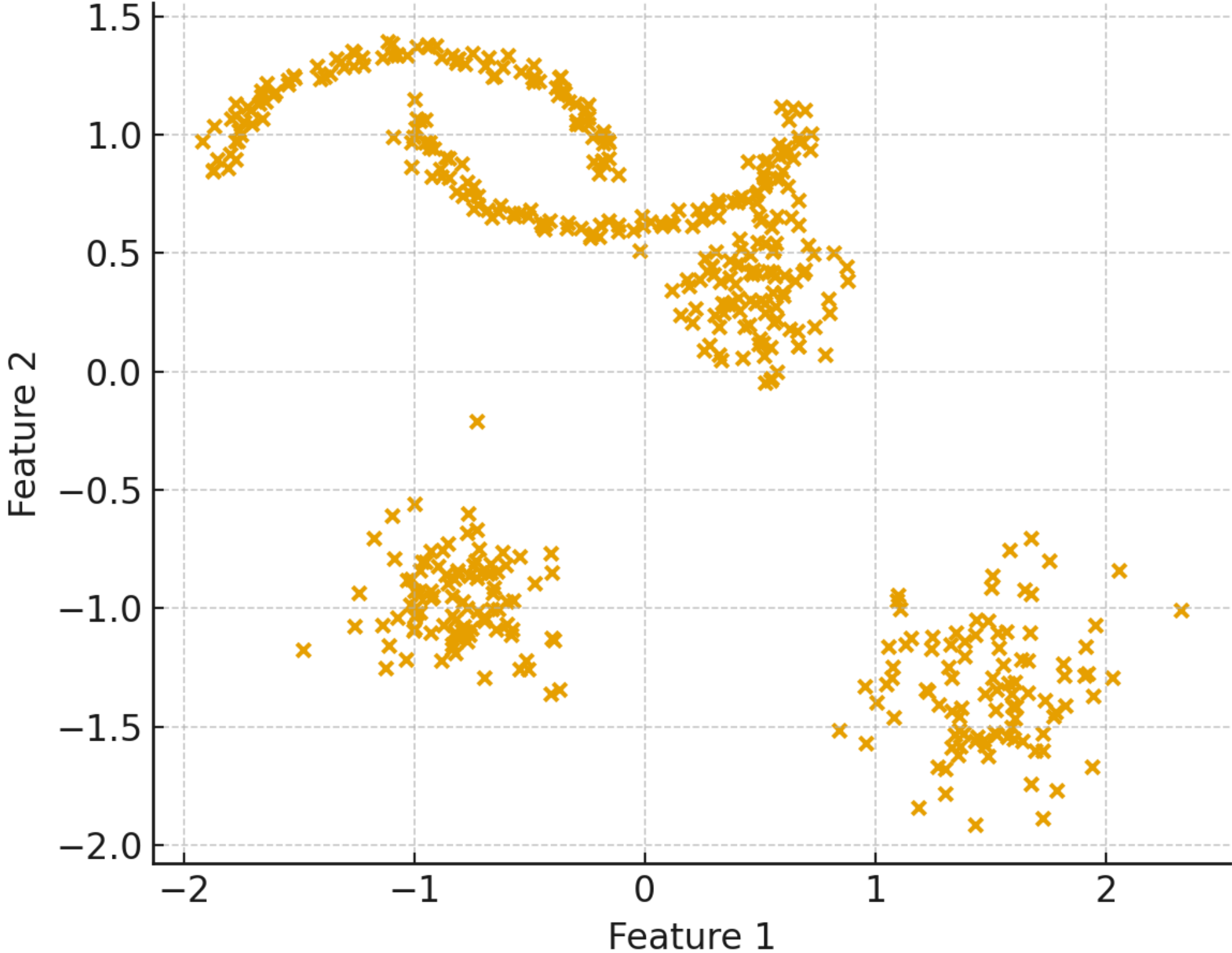Repeat until centroid movement $\approx 0$.

### Complexity

O(nkd,i) for (n) samples, (k) clusters, (d) dims, (i) iterations.
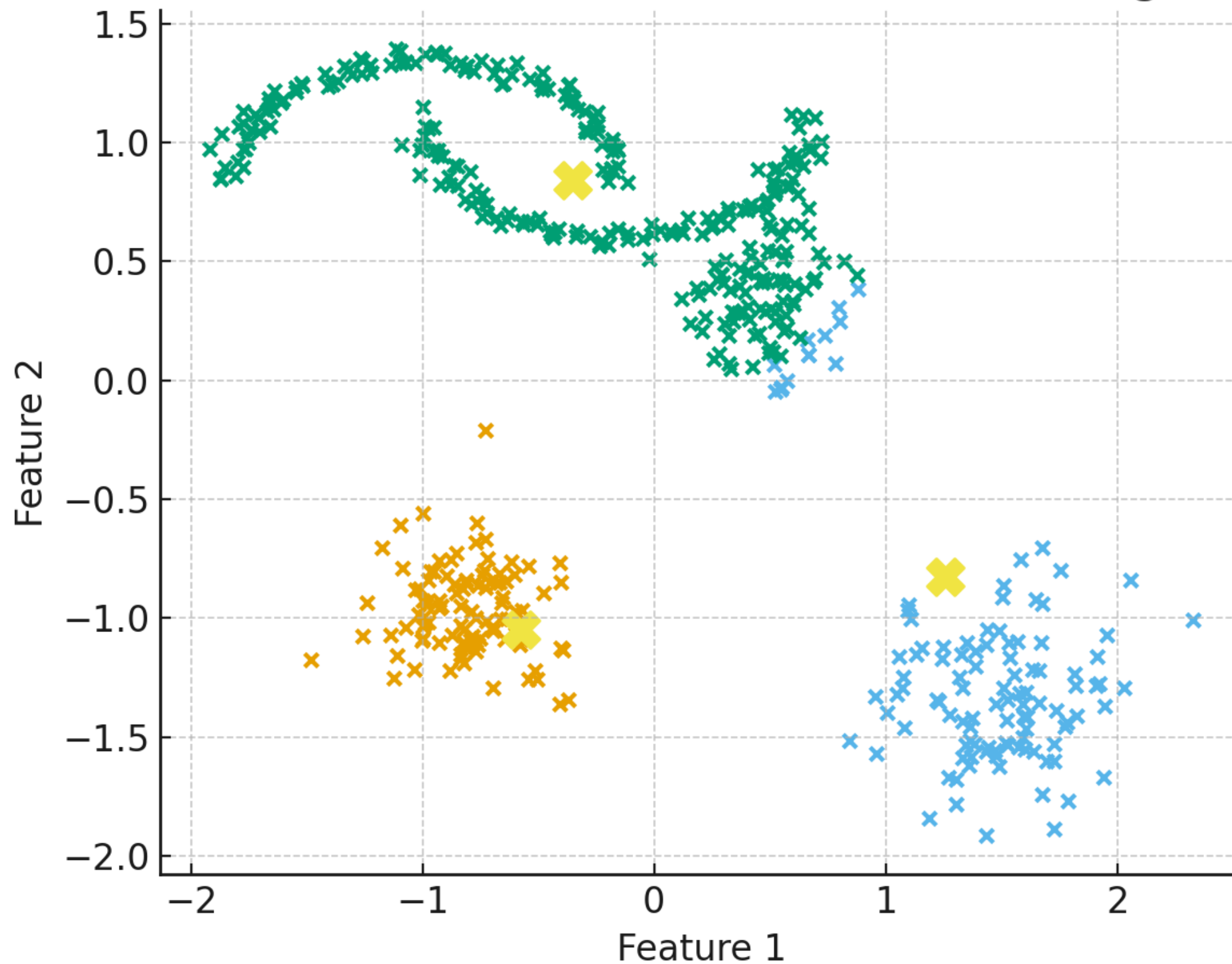
### Strengths / Limitations

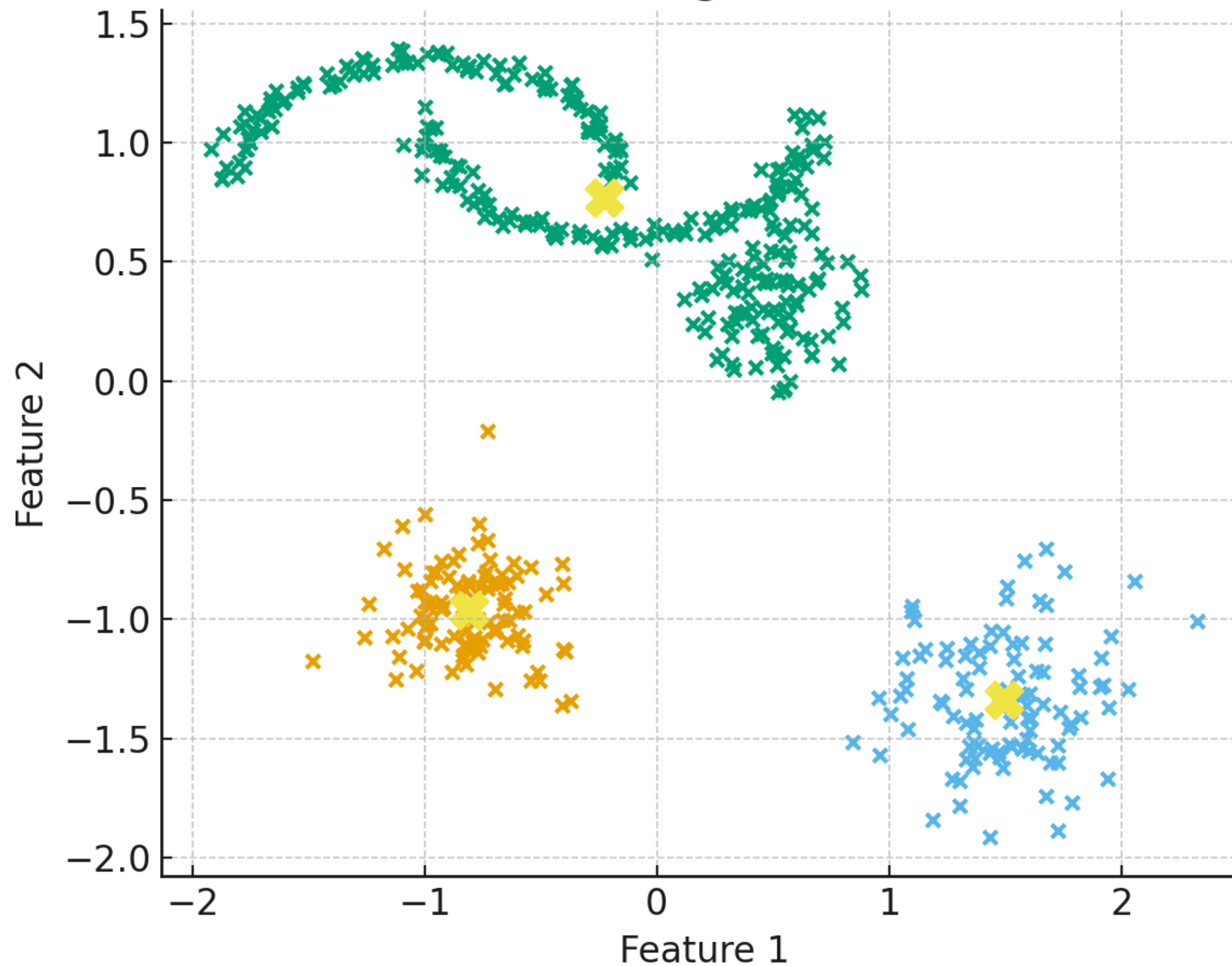fast and simple; assumes spherical clusters, sensitive to initialization and chosen (k).

K-Means Convergence — Raw Unlabeled Data

K-Means — Iteration 1 (initial centroids & assignments)

# Hierarchical Clustering & DBSCAN

## Hierarchical clustering (dendrograms) and density-based clustering (DBSCAN)
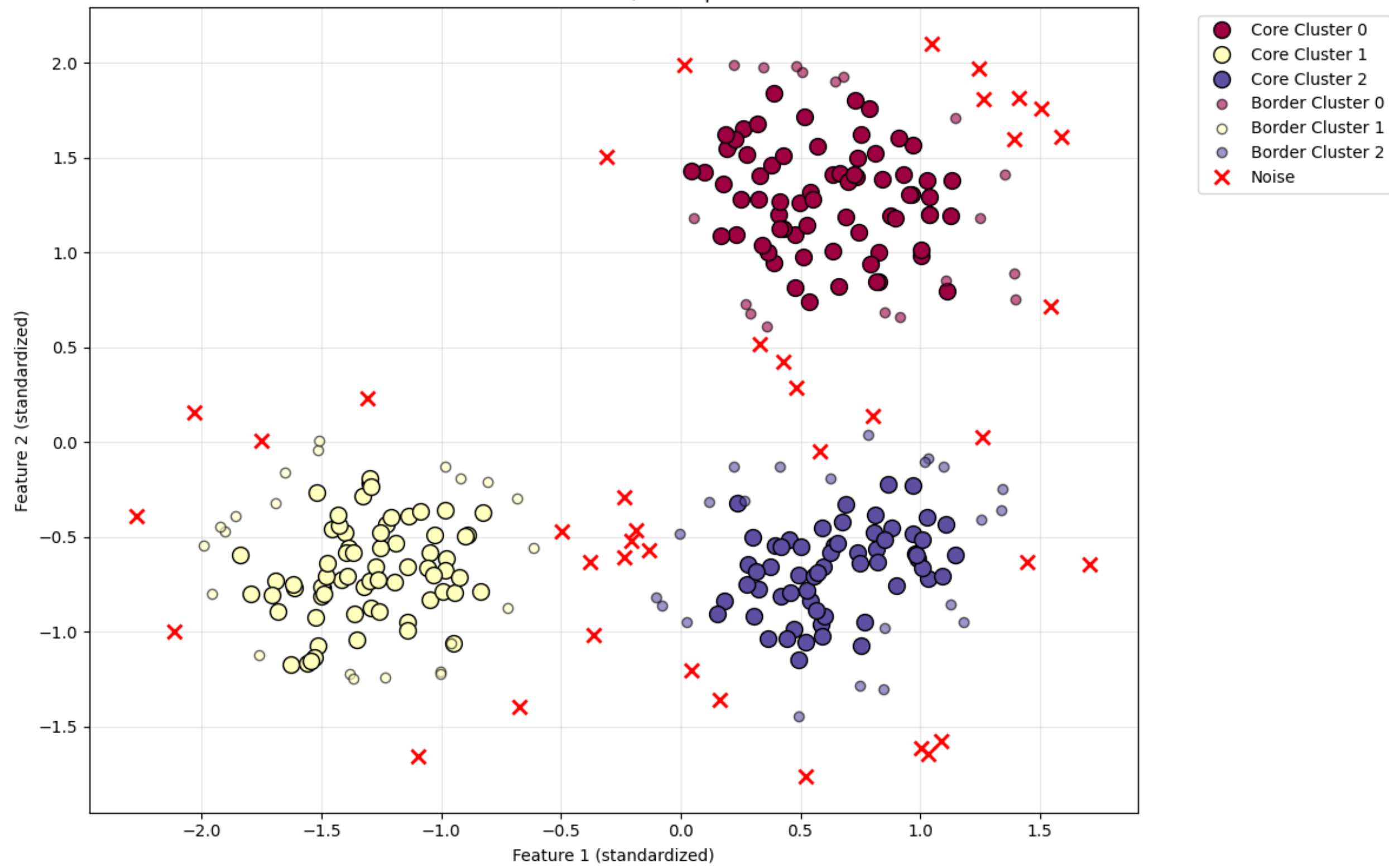
### (1) Hierarchical (tree-based)

- **Concept:** produce a nested cluster tree (dendrogram); no fixed (k) required.
- **Agglomerative (bottom-up):** start with singleton clusters and merge by linkage:
  - **Linkages:** Ward (minimize increase in WCSS), complete, average, single.
- **Divisive (top-down):** recursively split the dataset.
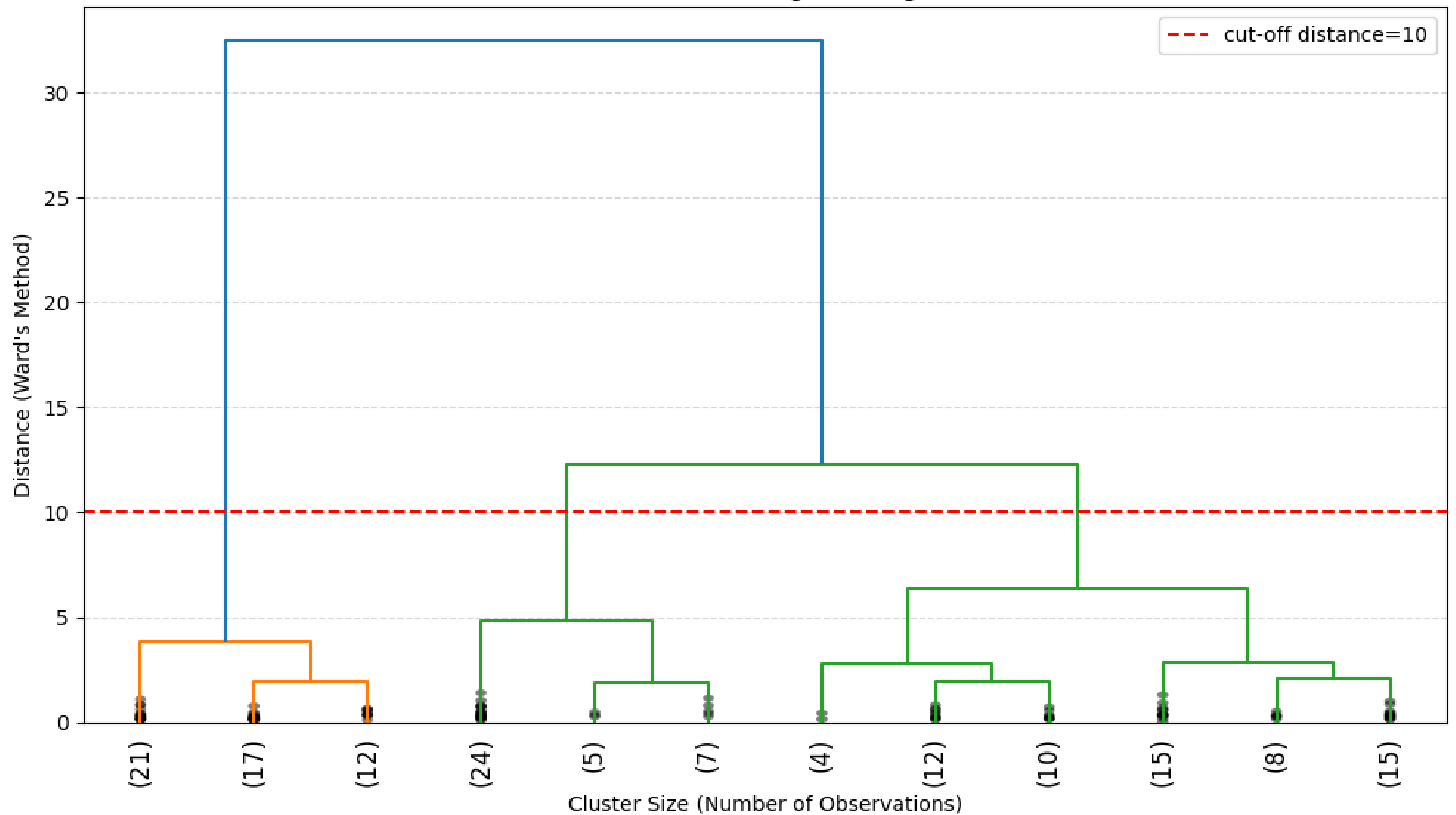- **Output:** dendrogram; choose cut height to extract clusters.

### (2) DBSCAN (density-based)

- **Parameters:** radius $\varepsilon$ and $min\_samples$.
- **Point types:**
  - **Core:** $|\mathcal{N}_\varepsilon(\mathbf{x})| \geq \text{min\_samples}$.
  - **Border:** $\varepsilon - neighborhood$ of a core but not a core.
  - **Noise:** neither core nor border.
- **Mechanics:** expand clusters from core points by density connectivity — finds arbitrary-shaped clusters and flags outliers.
- **Pros / Cons:** robust to noise; no (k) required; sensitive to $\varepsilon$ and varying density.

DBSCAN Clustering
Estimated clusters: 3, Noise points: 39

Truncated Hierarchical Clustering Dendrogram (Iris Dataset)

# Internal Clustering Evaluation Metrics

To assess the quality of clustering results objectively, internal metrics evaluate the structure inherent in the data without requiring external ground-truth labels.

**1**

## Silhouette Coefficient

Measures how similar an object is to its own cluster (cohesion) compared to other clusters (separation), based on pairwise distances.

**2**

## Davies-Bouldin Index

Evaluates the average "similarity" between each cluster and its most similar one, considering the dispersion within clusters and separation between their centroids.

**3**

## Calinski-Harabasz (CH) Index

Calculates the ratio of between-clusters dispersion to within-cluster dispersion across all clusters, analogous to an F-statistic in ANOVA.

# Internal clustering evaluation

## Internal metrics (no ground truth)



### ① Silhouette coefficient

Per-sample:

$$s(i) = \frac{b(i) - a(i)}{\max a(i), b(i)},$$

where $a(i)$ = mean intra-cluster distance, $b(i)$ = mean nearest-cluster distance.

**Range:** (-1) (misplaced) to (+1) (well separated). Use the mean $s(i)$ as a global score.
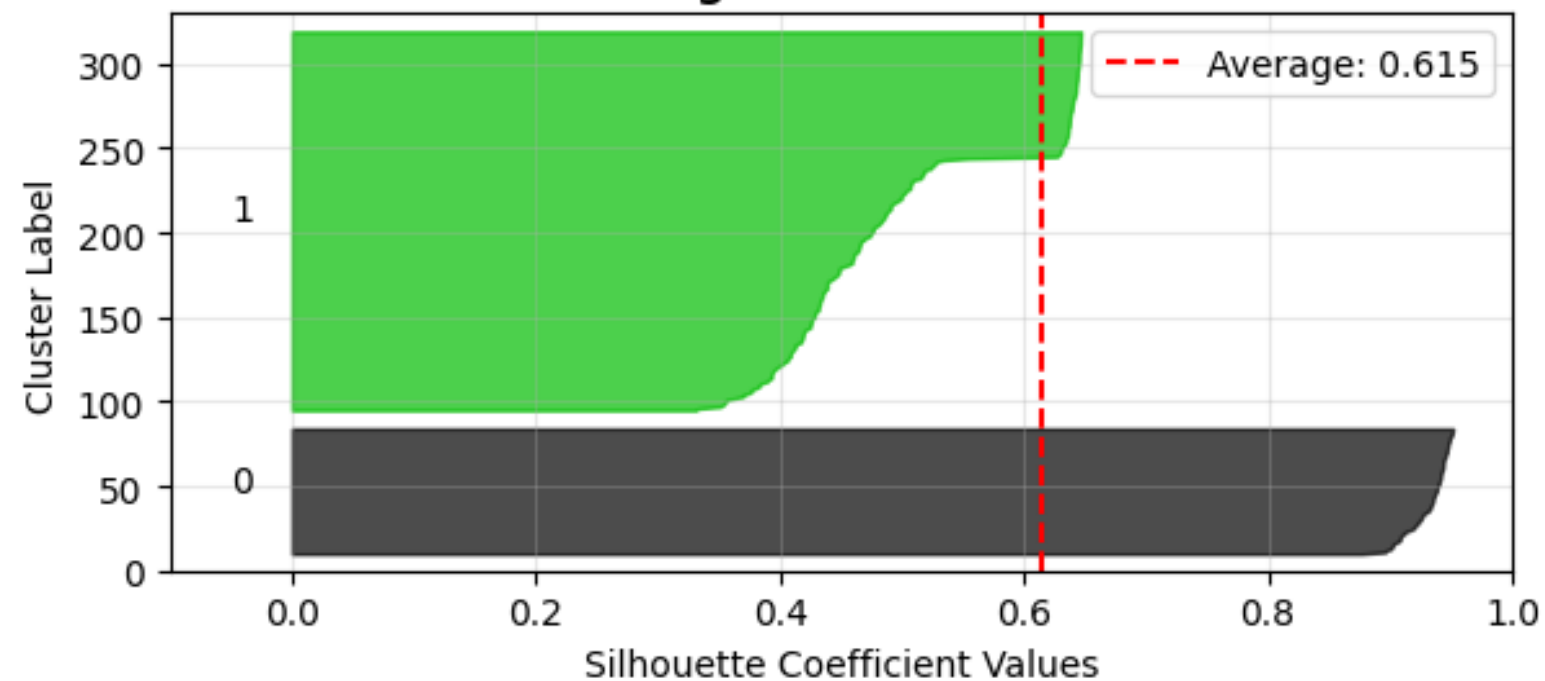
### ② Davies–Bouldin index (DB)

For (k) clusters:

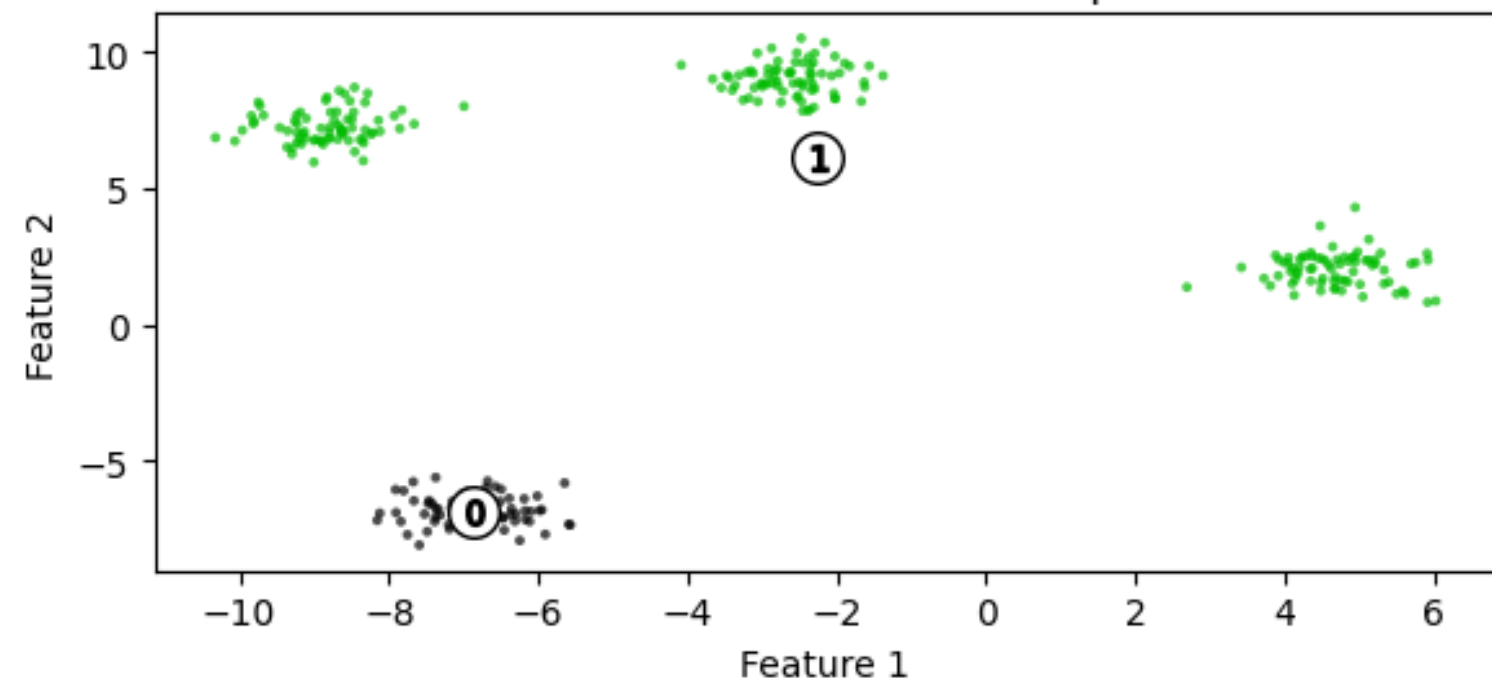$$DB = \frac{1}{k} \sum_{i=1}^{k} \max_{j \neq i} \frac{\sigma_i + \sigma_j}{d(c_i, c_j)},$$

where $\sigma_i = average\ distance$ of cluster (i) points to centroid $c_i$, $d(\cdot,\cdot) = centroid\ distance$.

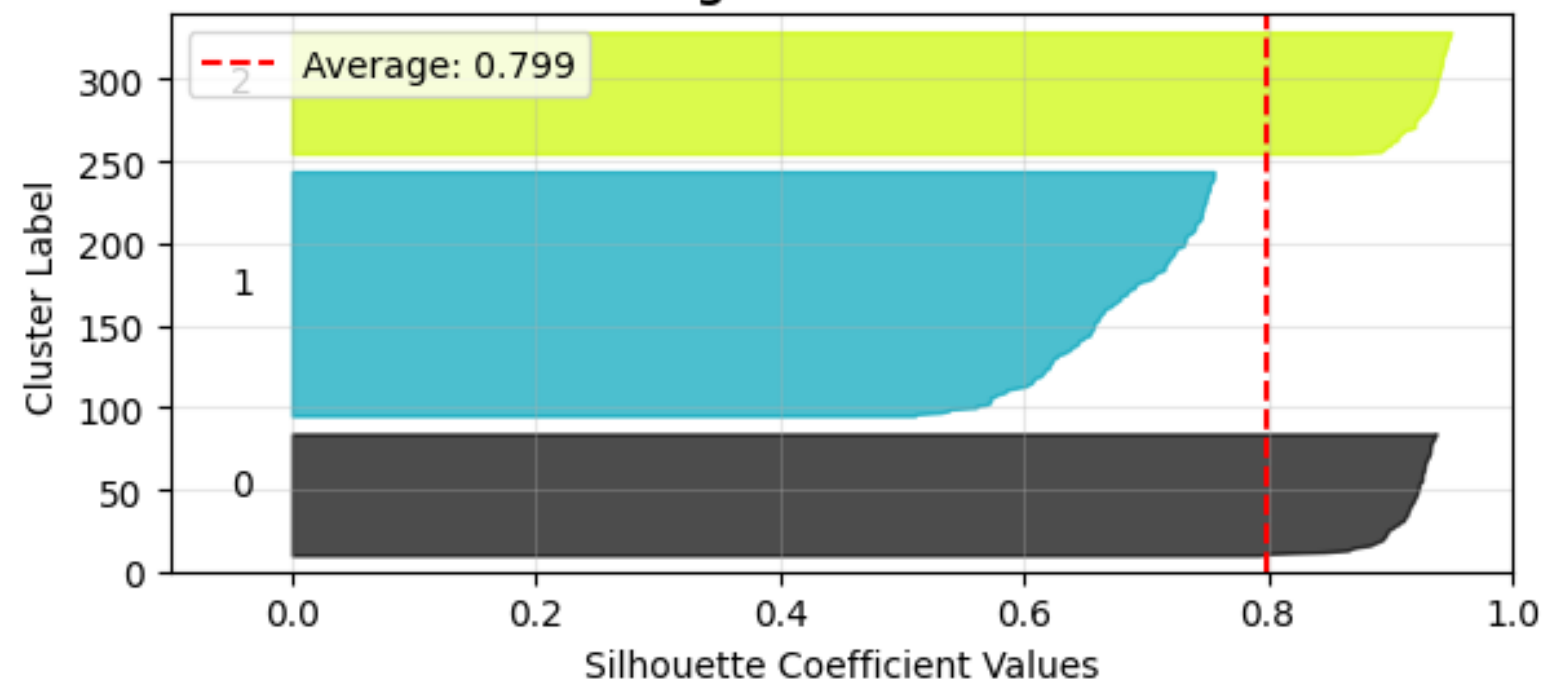**Interpretation:** lower DB = better (compact & well separated).

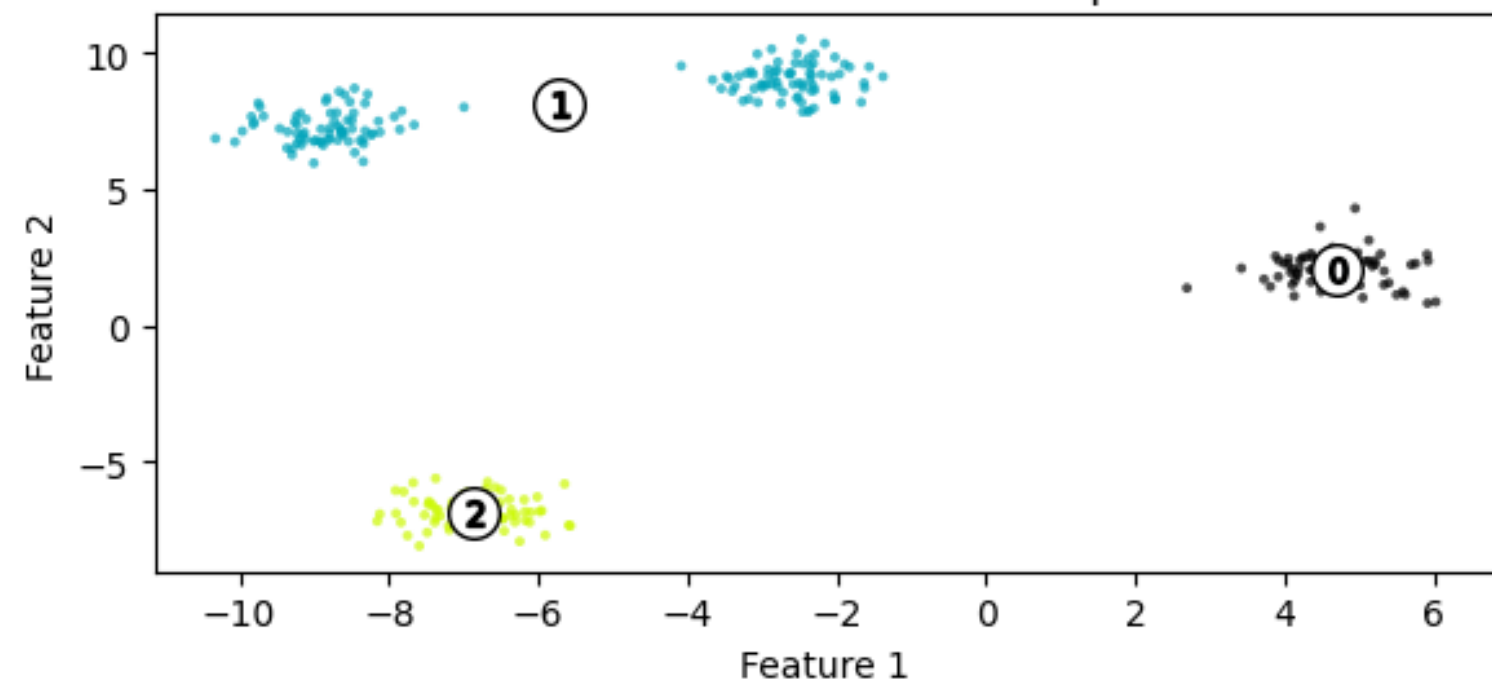**Silhouette Analysis for k = 2**
**Avg Score: 0.615**

**Data Clustered into 2 Groups**

**Silhouette Analysis for k = 3**
**Avg Score: 0.799**

**Data Clustered into 3 Groups**

**Silhouette Analysis for k = 4**
**Avg Score: 0.876**

Data Clustered into 4 Groups
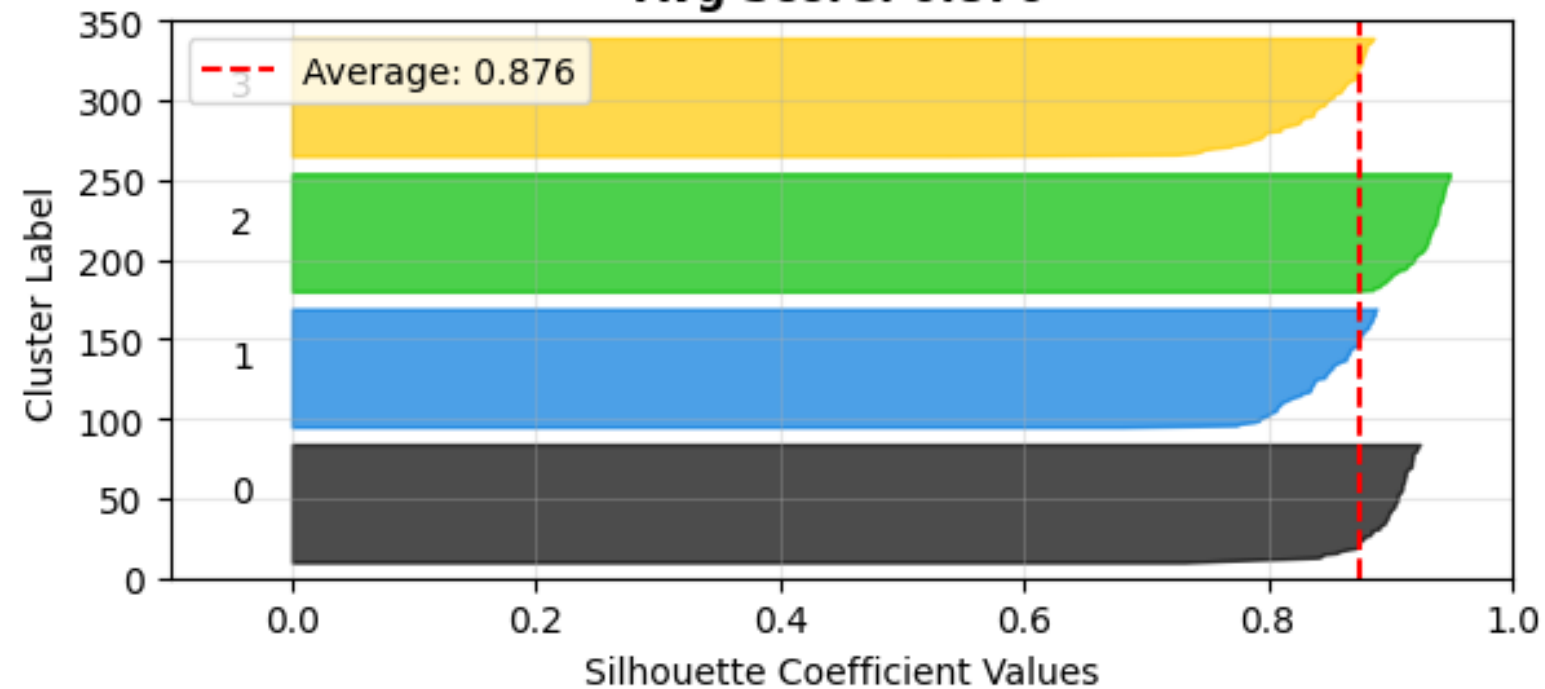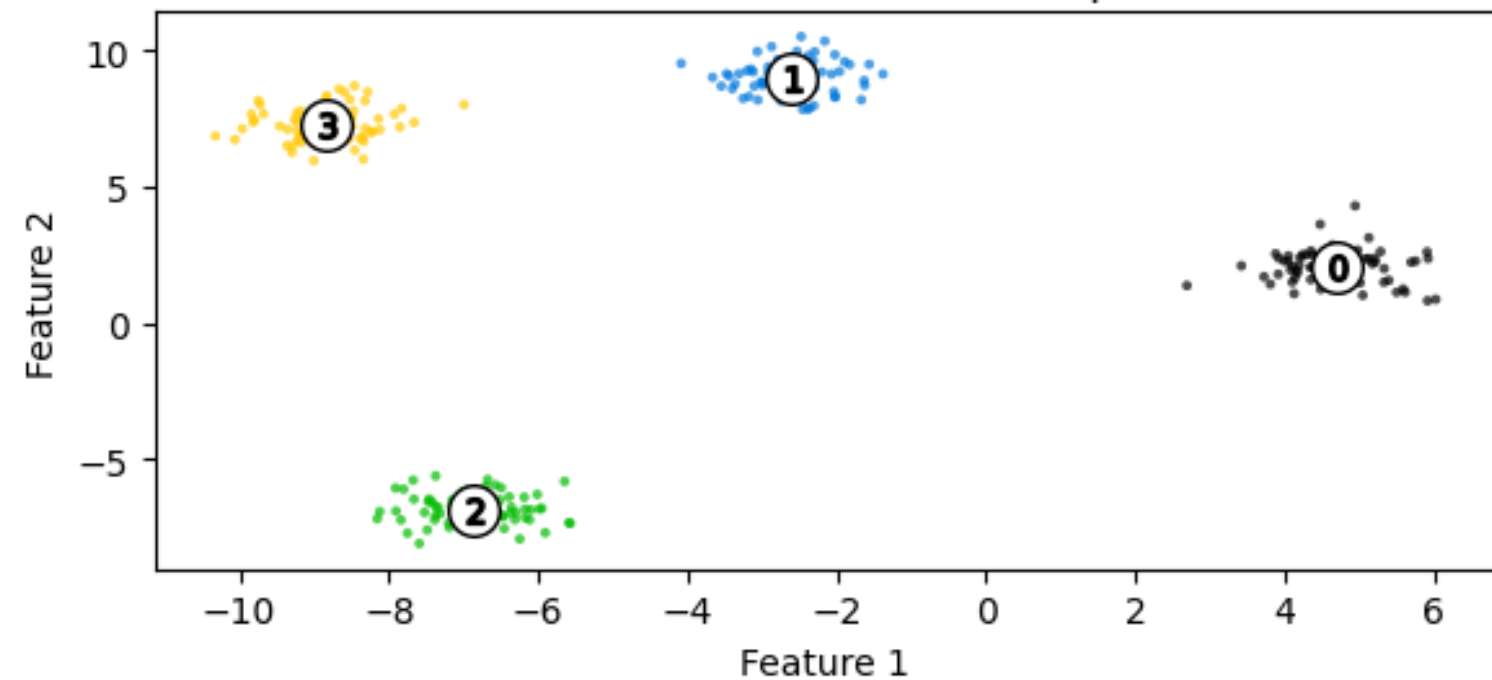
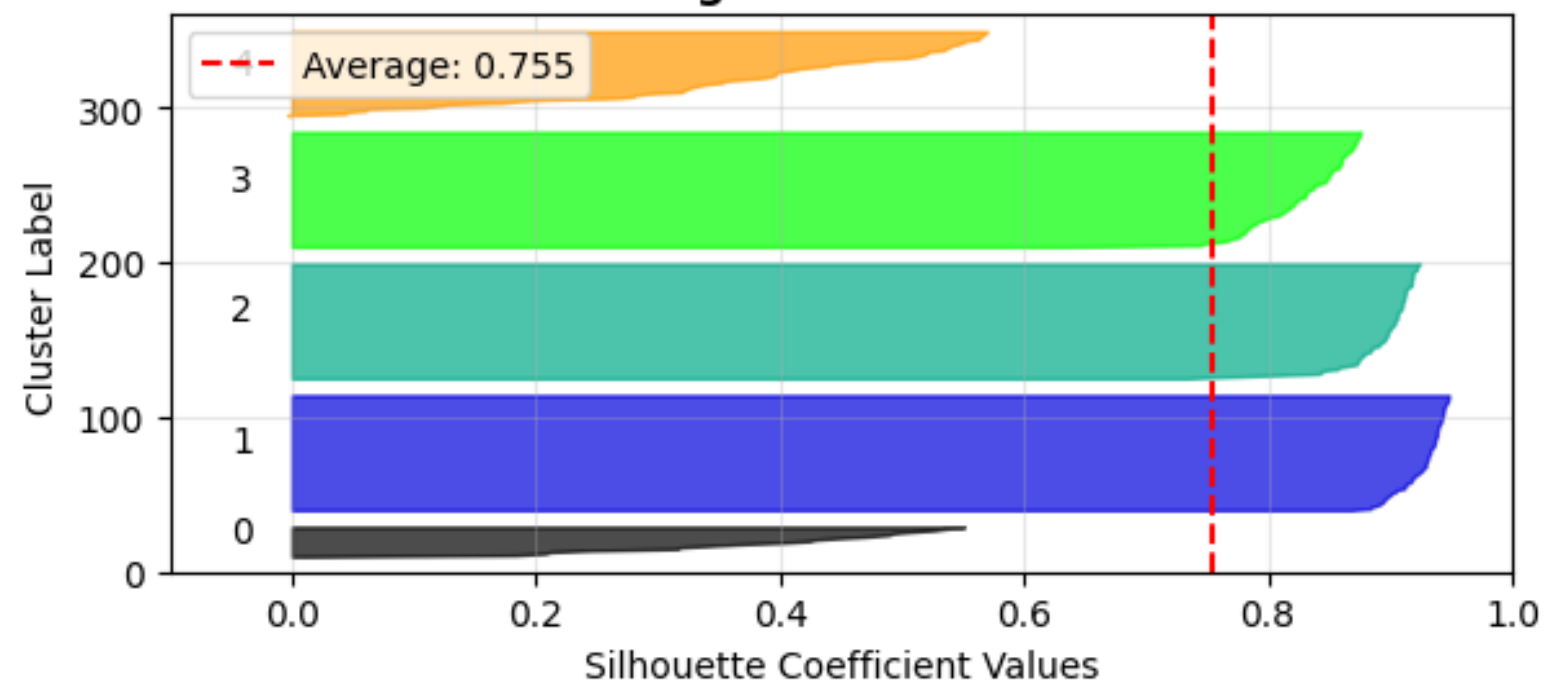**Silhouette Analysis for k = 5**
**Avg Score: 0.755**

Data Clustered into 5 Groups

Silhouette Analysis for k = 6
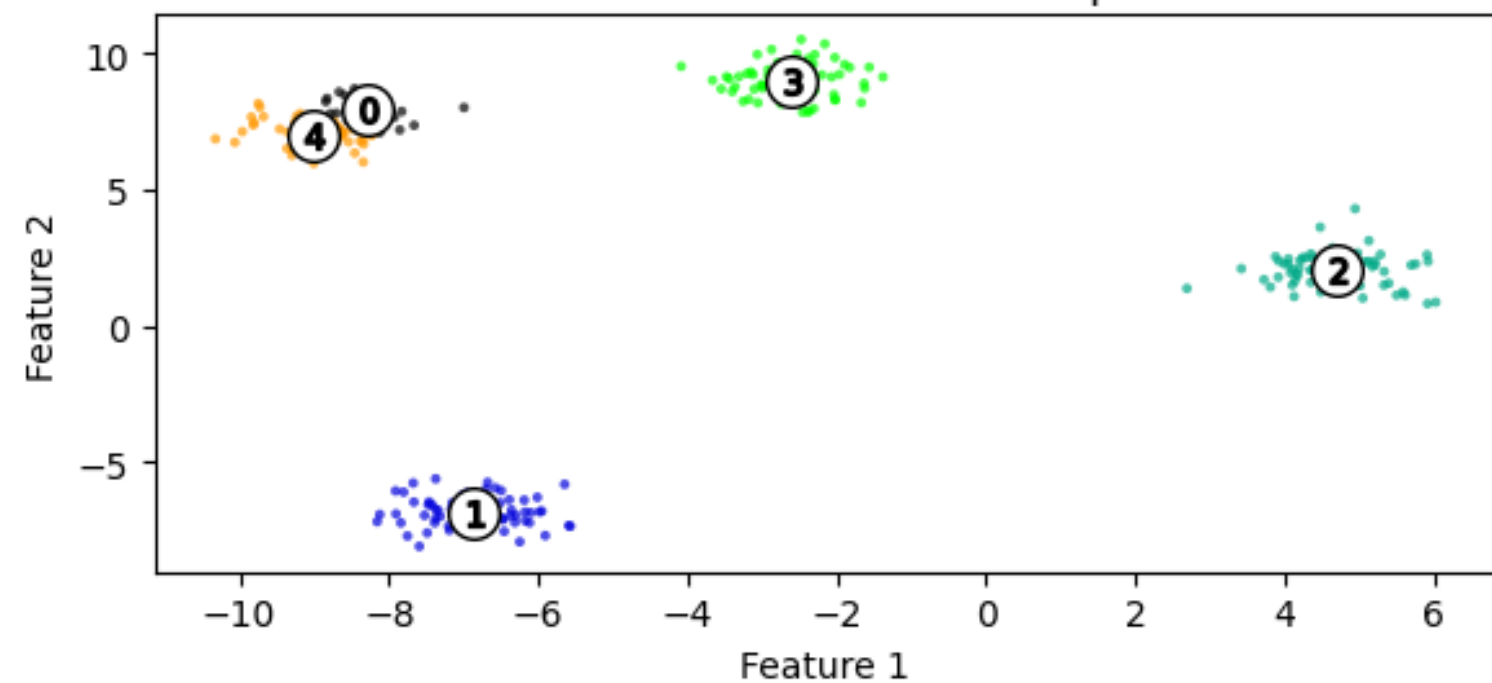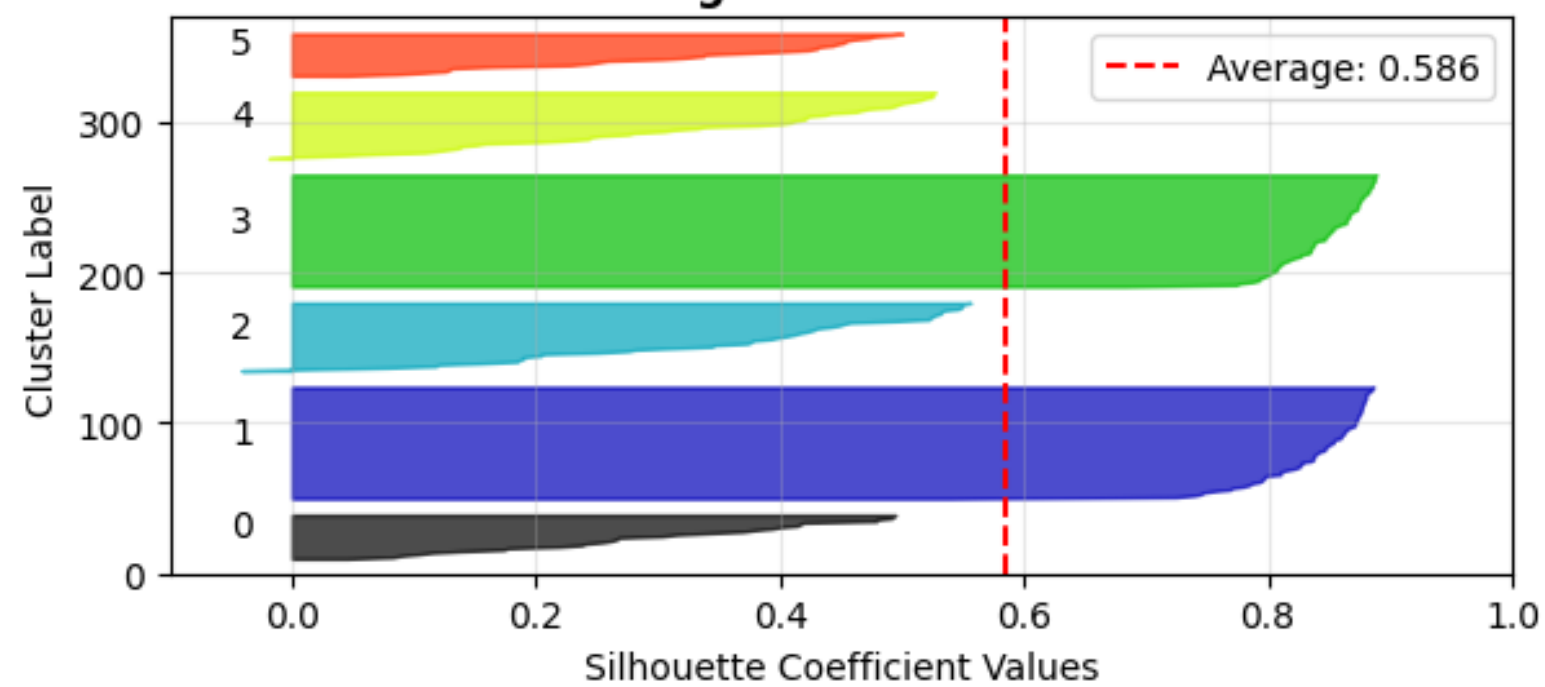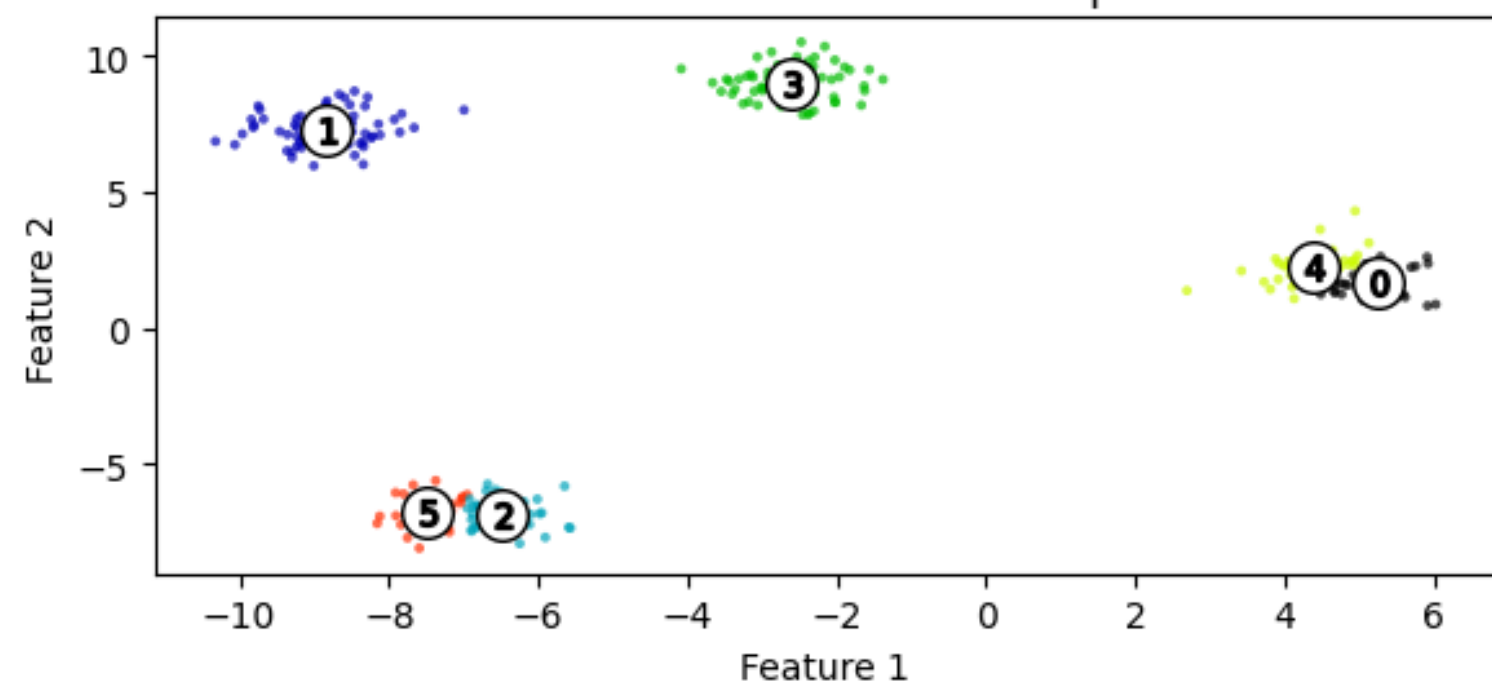Avg Score: 0.586
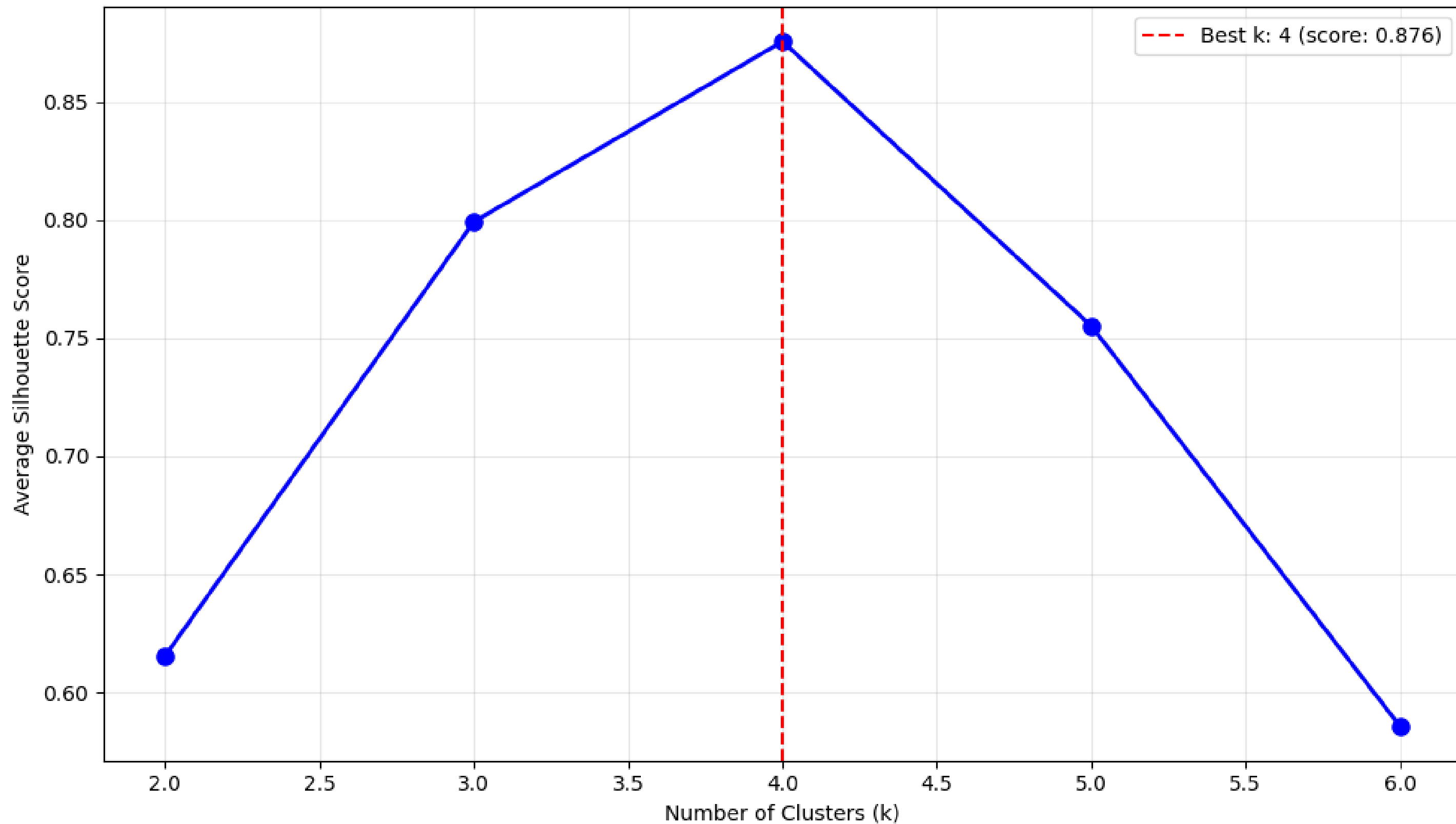
Data Clustered into 6 Groups

# Additional metrics & practical comparison

## Calinski–Harabasz and metric comparison notes

### Calinski–Harabasz (CH)

Definition:

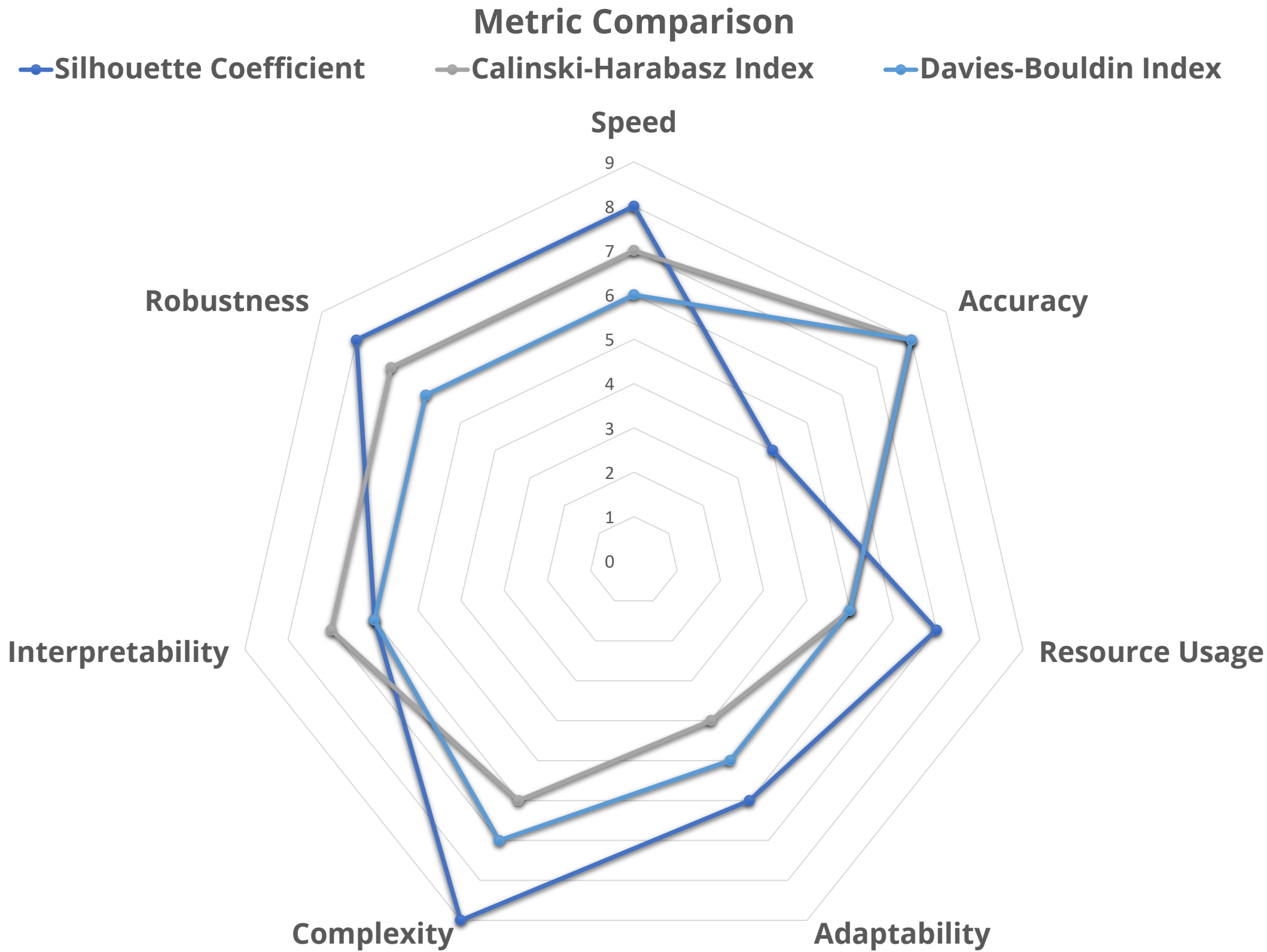$$CH = \frac{\text{trace}(B_k)/(k-1)}{\text{trace}(W_k)/(n-k)}$$

where $B_k$ = between-cluster dispersion matrix, $W_k$ = within-cluster dispersion matrix.

**Interpretation:** larger CH = better separated, compact clusters.

### Metric goals (practical checklist)

- **Maximize:** Silhouette, CH (higher → better cohesion/separation).
- **Minimize:** Davies–Bouldin (lower → better).

> 🗒 **Caveat:** metrics have different numeric scales—compare relative rankings or normalize before aggregation.
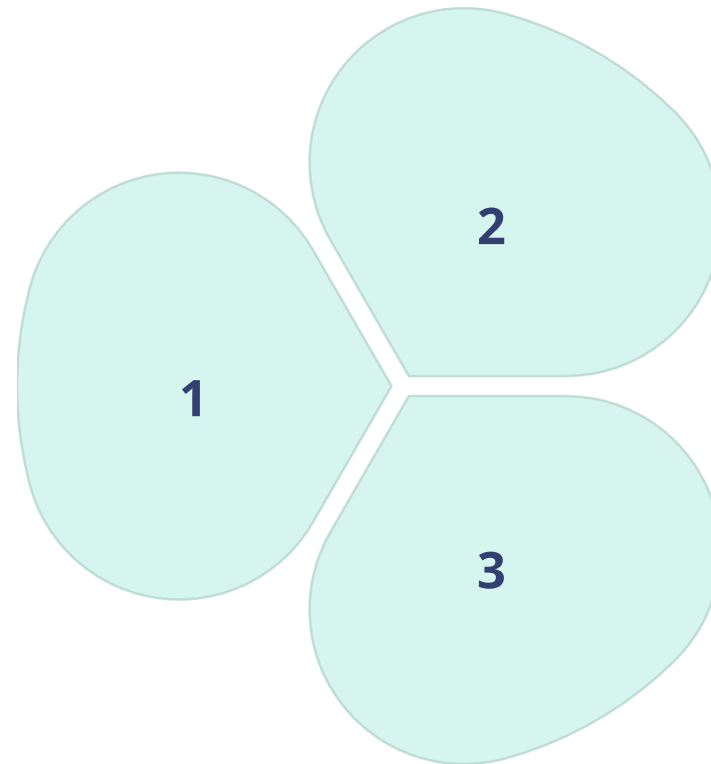
Metric Comparison

Silhouette Coefficient    Calinski-Harabasz Index    Davies-Bouldin Index

# Applications & concise conclusion

## Key applications and takeaways

### Applications

**Anomaly detection**

treat DBSCAN noise or low-density / far-from-centroid points as anomalies.

**2**

**Recommendation systems**

user segmentation (cluster users by behavior) and item clustering (recommend within same cluster).

**1**

**3**

**Feature engineering & EDA**

clustering to derive categorical features or reduce label sparsity.

### Conclusion

Clustering converts unlabelled data into structure—choose algorithm and metric based on shape, density, interpretability, and downstream needs.