

Lesson 29 Summary: K-Means Clustering

Data Science with Python – Key Concepts

Data Science Program

K-Means Clustering

Algorithm

1. Init centroids
2. Assign points
3. Update centroids

Choosing K

Elbow method
Silhouette score
Domain knowledge

Limitations

Spherical clusters
Sensitive to init
Need to specify K

Initialization

K-means++ (default): smart init
n_init=10: run multiple times

Finance Use Cases

Customer segmentation
Stock clustering | Portfolio groups

`KMeans(n_clusters=K).fit(X) | model.labels_ | model.cluster_centers_`

K-Means is the most popular clustering algorithm

Iterative optimization:

- **Step 1:** Initialize K centroids (randomly or K-means++)
- **Step 2:** Assign each point to nearest centroid
- **Step 3:** Update centroids as cluster means
- **Repeat:** Steps 2-3 until convergence

Algorithm converges when assignments stop changing

Methods to find optimal clusters:

- **Elbow method:** Plot inertia vs K, find bend
- **Silhouette score:** Measures cluster cohesion
- **Domain knowledge:** Business requirements

There is no single correct K; use multiple methods

K-means++ (default):

- **Problem:** Random init can give poor results
- **Solution:** Smart initialization spreads centroids
- **n_init:** Run multiple times, keep best

K-means++ is much better than random initialization

Critical for K-means:

- **Why:** Distance-based algorithm
- **Without scaling:** Large-scale features dominate
- **Solution:** StandardScaler before clustering

Always scale features before K-means

When K-means struggles:

- **Non-spherical:** Assumes circular clusters
- **Varying sizes:** All clusters treated equally
- **Outliers:** Can distort centroids

Consider DBSCAN or hierarchical for complex shapes

Evaluation metrics:

- **Inertia:** Sum of squared distances to centroids
- **Silhouette:** Cohesion vs separation (-1 to 1)
- **Higher silhouette:** Better-defined clusters

Silhouette close to 1 indicates well-separated clusters

Clustering use cases:

- **Customer segmentation:** Group by behavior
- **Stock clustering:** Similar price patterns
- **Regime detection:** Market states

Clustering reveals hidden structure in financial data

Essential Commands:

Task	Code
Fit model	<code>KMeans(n_clusters=K).fit(X)</code>
Get labels	<code>model.labels_</code>
Get centroids	<code>model.cluster_centers_</code>
Inertia	<code>model.inertia_</code>
Silhouette	<code>silhouette_score(X, labels)</code>

K-means is fast, simple, and effective for spherical clusters