

# Lesson 29: K-Means Clustering

Data Science with Python – BSc Course

Data Science Program

BSc Course

45 Minutes

# The End of Labels

Supervised learning needed labels for every example.

- L21–L28: every row had a target – price, class, fraud/not-fraud
- But real-world data is often **unlabeled**
- Can we still learn something useful?



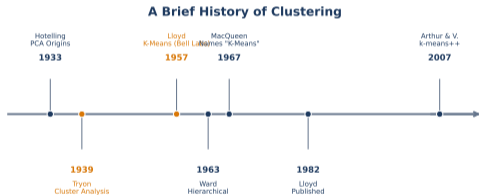
---

Welcome to unsupervised learning – discovering structure without being told what to look for

# Brief History of Clustering

## Clustering has a rich history across many disciplines

- 1939: Tryon applied cluster analysis in psychology research
- 1957: Lloyd invented K-Means at Bell Labs (published 1982)
- 1967: MacQueen coined the name “K-Means”



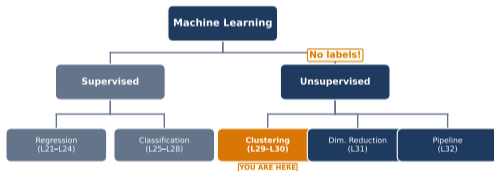
---

K-Means was invented at Bell Labs in 1957 but not formally published until 25 years later

# Where Unsupervised Fits

## ML has two main paradigms

- **Supervised** (L21–L28): predict a known target
- **Unsupervised** (L29–L32): discover hidden structure
- Today: clustering – grouping similar observations



---

Unsupervised learning finds patterns without any labels – clustering is its most common task

# Learning Objectives

**The Problem:** We have 500 stocks with dozens of features. How do we group similar stocks without predefined categories?

**After this lesson, you will be able to:**

1. Explain the K-Means algorithm step by step
2. Choose optimal K using the elbow method and silhouette score
3. Interpret cluster centroids to understand group characteristics
4. Apply K-Means to segment financial assets by behavior

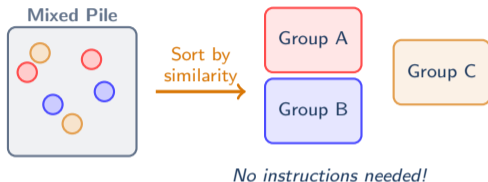
**Prerequisites:** L13 (descriptive statistics), L21 (distance concepts)

---

**Finance application:** stock segmentation, customer clustering, market regime detection

# What is Clustering?

Sorting without instructions – like organizing laundry



- **Goal:** group similar items together, separate dissimilar ones
- **Key insight:** the algorithm discovers groups on its own
- **Finance:** group stocks by behavior, not by sector label

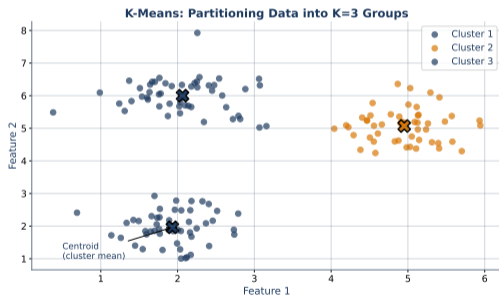
---

Clustering = grouping by similarity. No one tells the algorithm what the groups should be.

# K-Means: Core Idea

## Partition data into exactly K groups

- You choose K (number of clusters) in advance
- Each point belongs to the cluster with the nearest center
- Centers (centroids) are updated iteratively until stable



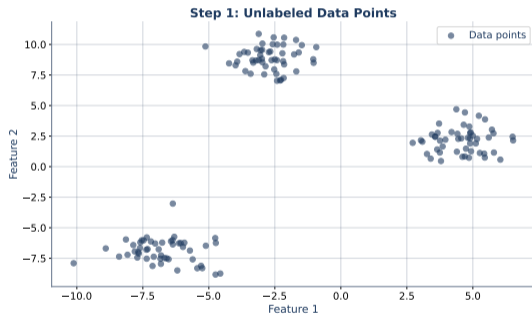
---

K-Means finds compact, spherical clusters by minimizing within-cluster variance (inertia)

# Step 1: Raw Data

## Starting point – unlabeled observations

- No colors, no groups – just scattered points
- Goal: discover natural structure hidden in the data



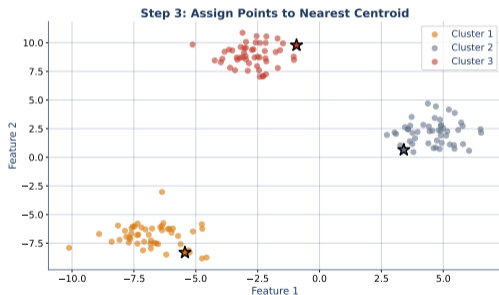
---

Before clustering: no structure visible. The algorithm will find groups automatically.

# Steps 2–3: Initialize & Assign

Place  $K$  centroids, then assign every point

- **Step 2:** Place  $K$  initial centroids (randomly or via k-means++)
- **Step 3:** Assign each point to its nearest centroid
- This creates a Voronoi partition of the feature space



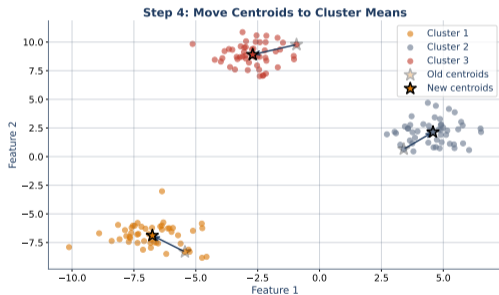
---

Each point is assigned to the cluster whose centroid is closest (Euclidean distance)

# Step 4: Recompute & Repeat

## Move centroids to cluster means, then reassign

- New centroid = mean position of all points in that cluster
- Repeat Steps 3–4 until assignments stop changing
- Convergence is guaranteed (finite iterations)



---

K-Means converges when no point changes cluster. Typically 10–20 iterations.

# How Many Groups?

**K-Means requires you to choose K in advance**

This is the central challenge:

- Too few clusters → merge distinct groups
- Too many clusters → split natural groups apart
- No single “correct” answer – it depends on the question

**Two complementary methods:**

1. **Elbow method** – plot inertia vs K, look for the bend
2. **Silhouette score** – measure how well-separated clusters are

**Finance example:**

Are there 3 market regimes (bull/bear/sideways) or 5 (adding crash and recovery)?

The data can support multiple valid answers.

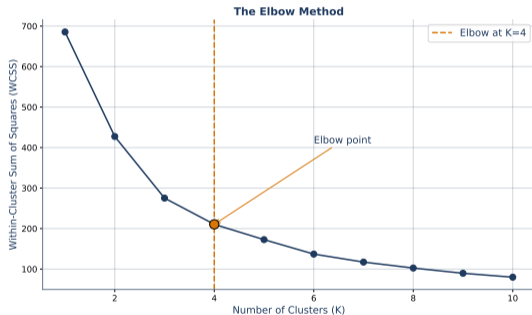
---

**Choosing K is the hardest part of K-Means. Use elbow + silhouette + domain knowledge.**

# The Elbow Method

Plot inertia vs K – look for diminishing returns

- Inertia always decreases as K increases
- The “elbow” = where adding more clusters stops helping much



The elbow is subjective – combine with domain knowledge and silhouette score

# Checkpoint: Test Your Understanding

**Q1:** What does inertia measure?

**Q2:** Why must we standardize before K-Means?

**Q3:** If the elbow plot has no clear bend, what does that suggest about the data?

## Quick answers:

1. Inertia = total squared distance from points to their cluster centers
2. Features on different scales would dominate distance calculations
3. The data may not have well-defined cluster structure

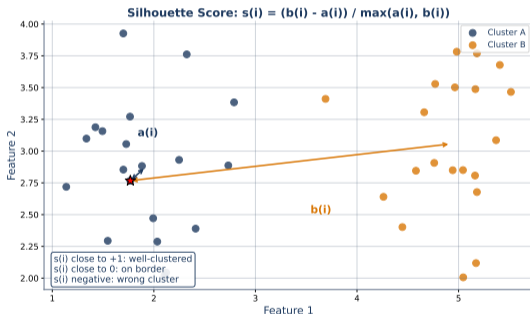
---

Pause here – make sure you can answer all three before continuing

# Silhouette Score

Is each point in the right cluster?

- Range:  $-1$  to  $+1$ . Higher = better-separated clusters
- Compares distance to own cluster vs nearest other cluster
- $s > 0.5$ : reasonable structure;  $s < 0$ : wrong cluster

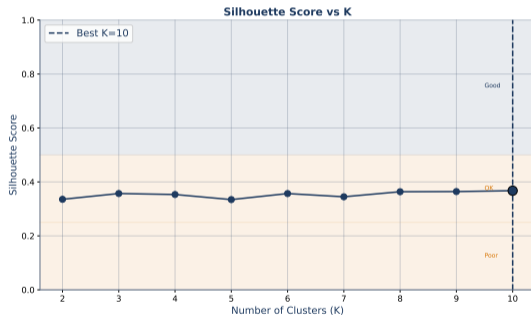


Silhouette answers: “Am I closer to my own cluster or to the next-nearest one?”

# Silhouette Score vs K

## Choose K that maximizes the average silhouette

- Compute silhouette for  $K=2, 3, 4, \dots$
- Peak = best trade-off between cohesion and separation
- Combine with elbow method for robust choice

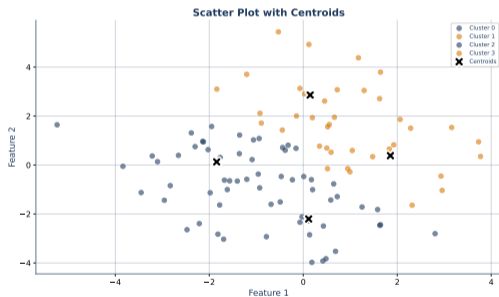


Silhouette peak + elbow bend at same K = strong evidence for that number of clusters

# Visualizing Clusters

## Scatter plot – the first sanity check

- Color points by cluster assignment
- Mark centroids with a distinct symbol
- For high-dimensional data: reduce to 2D with PCA first

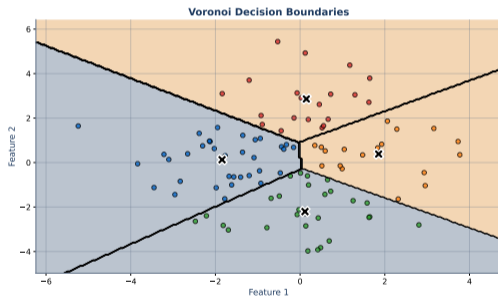


Always visualize clusters to check whether the groupings make intuitive sense

# Voronoi Boundaries

## Where one cluster ends and another begins

- Each region contains all points closest to that centroid
- Boundaries are straight lines (linear decision boundaries)
- Like territorial borders based on distance to capitals

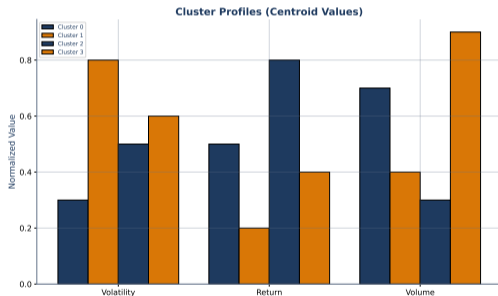


Voronoi partitions create linear boundaries – K-Means cannot capture non-convex shapes

# Interpreting Centroids

## Centroids tell you what each cluster “looks like”

- Centroid = average feature values for all points in that cluster
- Compare across clusters to understand what distinguishes them
- Use `scaler.inverse_transform()` to get original-scale values



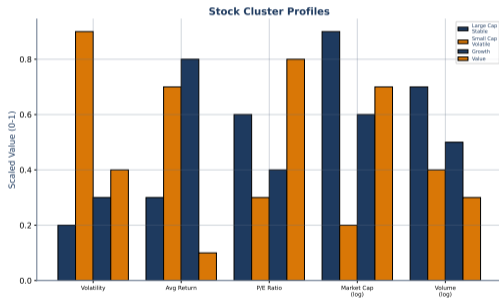
---

Bar chart of centroid features reveals what makes each cluster distinct

# Finance: Stock Cluster Profiles

## What does each stock group look like?

- Features: return, volatility, beta, market cap, P/E ratio
- Centroids reveal natural groupings beyond sector labels
- Name clusters: “Blue Chips”, “Growth”, “Speculative”

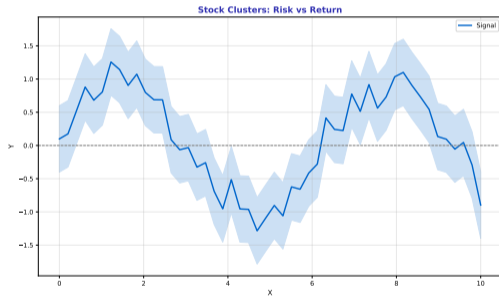


Cluster names come from interpreting centroids: high vol + low cap = “Speculative”

# Finance: Risk-Return Map

## Stocks clustered by volatility and return

- Each cluster occupies a different risk-return region
- Diversification strategy: pick one stock per cluster
- Clusters often cut across traditional sector boundaries



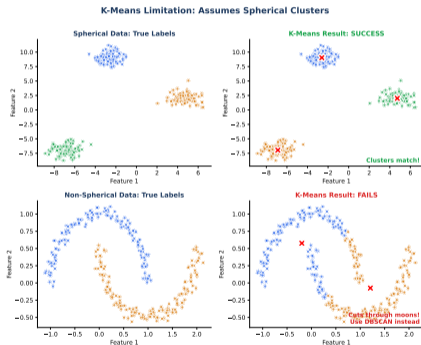
---

**Portfolio tip: select one stock per cluster for maximum diversification**

# K-Means Limitations

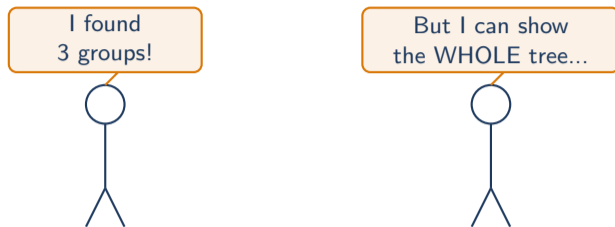
## When K-Means struggles

- Assumes spherical, equally-sized clusters
- Sensitive to outliers; cannot handle non-convex shapes



Works best for compact clusters. Hierarchical clustering (L30) handles complex shapes.

# What's Next: The Clustering Family



*Next: L30 Hierarchical Clustering*

**K-Means gives you flat groups. But what if groups have sub-groups?**

- Hierarchical clustering builds a tree of nested clusters
- No need to choose K in advance – the dendrogram shows all levels
- We will compare both methods in L30

---

**K-Means = fast, flat partitioning. Hierarchical = slower, richer structure. Both are essential.**

# Hands-On Exercise (25 min)

## Task: Cluster Stocks by Behavior

1. Calculate features: 1-year return, volatility, beta for 50 stocks
2. Standardize features using `StandardScaler`
3. Run K-Means with  $K = 2, 3, 4, 5$  – plot the elbow curve
4. Choose best  $K$  using elbow + silhouette score
5. Interpret centroids: what characterizes each cluster?

**Deliverable:** Elbow plot + scatter plot of clusters with labels

**Extension:** Compare clusters to GICS sectors – do behavioral clusters align with industry classifications?

---

**Remember:** standardize first, then cluster. Use `random_state=42` for reproducibility.

# Summary

**Problem solved:** Discovered natural groups in data without labels

## **Key takeaways:**

1. K-Means iterates: assign points to nearest centroid → recompute centroids → repeat
2. Choose K with the elbow method (inertia) and silhouette score
3. Always standardize features before clustering
4. Interpret centroids to name and understand clusters

## **Finance insight:**

Stock clusters often cut across sectors – behavioral groupings reveal diversification opportunities that sector labels miss.

**Next lesson:** L30 Hierarchical Clustering – clusters within clusters

---

**K-Means: fast, simple, effective. But choose K carefully and always standardize first.**