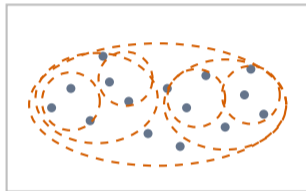


How many groups
are really there?



No labels. No ground truth. Just points.

What will you be able to do by the end?

By the end of this lecture you will be able to:

- **Analyze** why clustering output is meaningless without validation.
- **Apply** silhouette, elbow, and gap statistic to judge whether a clustering is real.
- **Evaluate** a clustering result by stability rather than a single silhouette number.
- **Compare** K-Means, Hierarchical, DBSCAN, and GMM on their different validation stories.

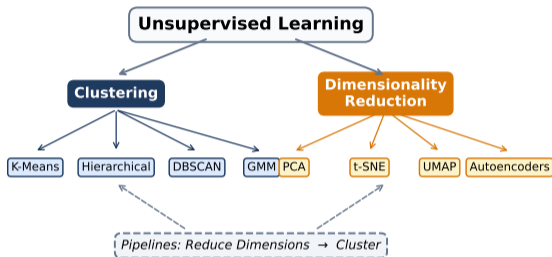
Retrieval: say one objective aloud in your own words.

What must you already know?

Before we begin: **try to recall** – what is the difference between supervised and unsupervised learning, in ONE sentence? Write it on paper. Do not look it up.

Retrieval: write your sentence before turning the page.

Why does unsupervised learning matter now?



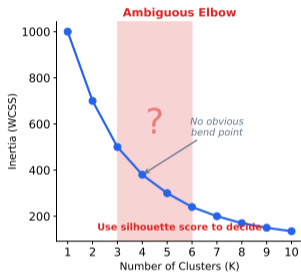
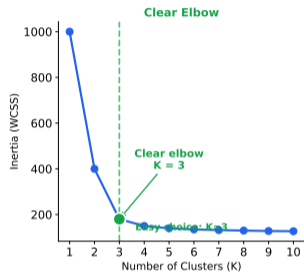
- Banks cluster customers into risk buckets **without** a predefined label.
- Hedge funds extract latent factors from 500 stocks with **no** “correct” factor list to check against.
- Card networks flag anomalies that **no** analyst pre-identified as fraud.

Unsupervised Learning Lecture

Retrieval: name one real-world task where you have data but no labels.

What is The Validation Paradox?

Elbow Method: Clear vs Ambiguous Cases



The Validation Paradox. You run unsupervised learning to discover groups you did not know existed. But you have no labels to verify your discovery is real.

Two values in conflict:

- **Discovery** – let the data speak; do not impose structure.
- **Validation** – to judge an output, you must impose assumptions (distance, K, density, covariance).

You cannot fully honor both.

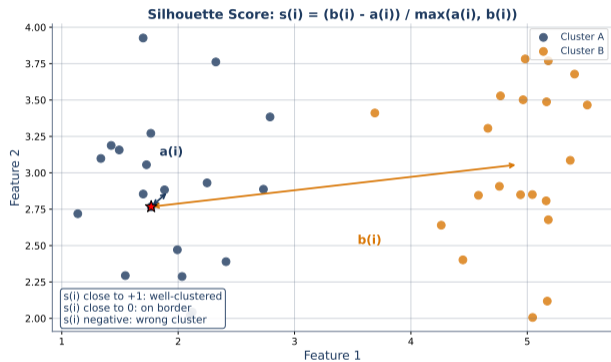
Retrieval: write the paradox's core tension in your own words.

Pause 1 – name the tension

Pause 30 seconds.

Write one sentence: *your boss hands you 1,000 customers and asks for 5 groups. How do you prove 5 is right and not 3 or 7?*

Where does the paradox come from, mathematically?

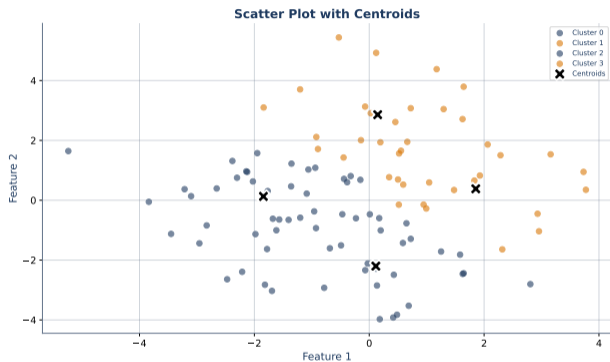


- Clustering has **no** loss-on-labels. Only internal scores: silhouette, within-cluster sum-of-squares, log-likelihood.
- **Silhouette**: an internal score in the range -1 to $+1$. Higher is better.
- All internal scores depend on a choice: distance metric, K , density threshold.

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

Retrieval: one sentence – what does a silhouette score measure?

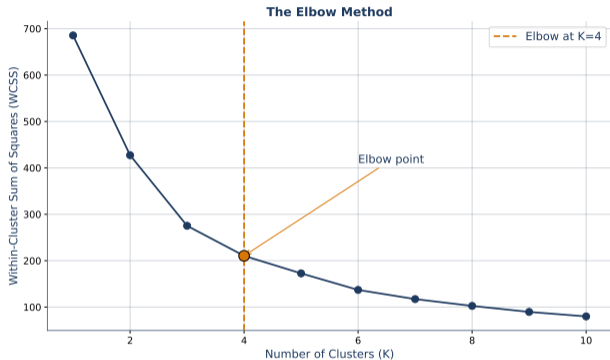
Can we cluster customers without knowing the segments?



- 1,000 customers, 50 behavior features, 0 labels.
- K-Means with $K = 5$ produces 5 clean buckets: “high-spend infrequent”, “mid-spend regular”, etc.
- Marketing asks: why 5 and not 3 or 7? How do you know these buckets are REAL?

Netflix re-clustered its 180M+ subscribers with deep autoencoders in 2023 after finding that K-Means segments were unstable quarter-to-quarter – a direct symptom of the validation paradox.

How can we validate a clustering without labels?



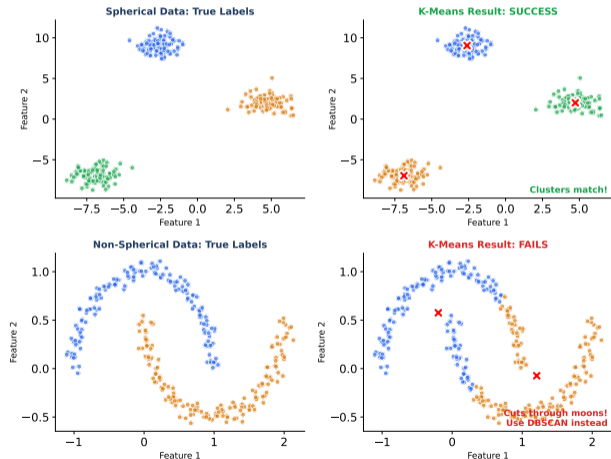
- A curve on K vs a loss metric, pick the inflection point.
- An average score over all points, pick the K with the highest value.
- Resample the data, cluster again, check if the assignments are stable.

$$WCSS(K) = \sum_{k=1}^K \sum_{x \in C_k} \|x - \mu_k\|^2$$

Retrieval: name two ways to pick K without labels.

What happens when two validation methods disagree?

K-Means Limitation: Assumes Spherical Clusters



- Some algorithms will return K groups by construction – whether or not K groups exist.
- **Disagreement:** silhouette says $K = 3$, elbow says $K = 5$; which do you trust?
- **Confirmation bias:** analysts pick the K that matches a story they already wanted to tell.

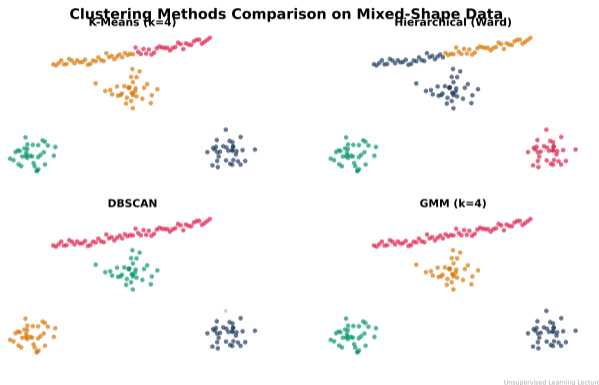
Pause 2 – resolve a disagreement

Pause 30 seconds.

Silhouette says $K = 3$. Elbow says $K = 5$.

Which do you believe, and why? What third test would break the tie?

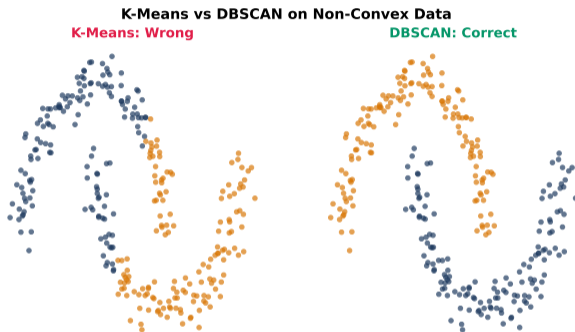
Which algorithm carries which validation story?



- **K-Means**: specify K up front; validate with silhouette/elbow.
- **Hierarchical**: no K needed; validate by dendrogram cut or cophenetic correlation.
- **DBSCAN**: no K ; returns “noise” points as a built-in discovery signal.
- **GMM**: soft assignments; validate with BIC/AIC – probabilistic, not geometric.

Retrieval: when would DBSCAN beat K-Means?

Who bears the cost of a “discovery” that was not real?



Unsupervised Learning Lecture

- **Retailers:** marketing to segments that do not exist wastes millions in campaign spend.
- **Clinicians:** mis-clustering patient phenotypes leads to wrong treatment arms in precision-medicine trials.
- **Traders/quants:** factor structures that are artifacts of PCA on too-few samples drive leveraged bets that unwind.

In 2022, a large US retailer retracted a segmentation study after internal data scientists showed the 7 “customer personas” produced by K-Means had near-zero silhouette AND did not persist across quarters. Budget impact: high 7 figures.

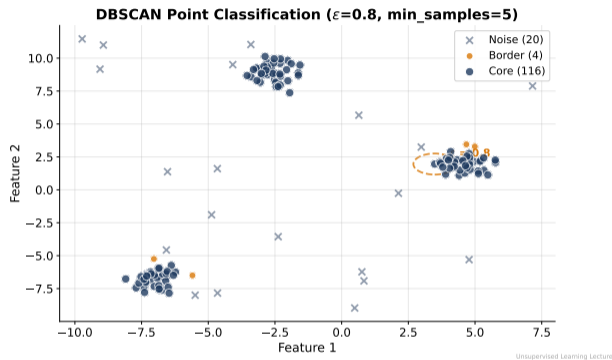
Pause 3 – interpret noise

Pause 30 seconds.

DBSCAN labels 20% of your customers as “noise” (not in any cluster).

Is that a good or bad result? What does it depend on?

So what counts as a real cluster?



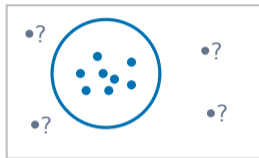
The Validation Paradox. Discovery without validation is just guessing.

- Persistence: the same grouping appears across algorithms, across bootstrap samples, across time windows.
- Story: the grouping corresponds to something a domain expert recognizes.
- The answer is always in a table, not a number.

Retrieval: in one sentence, when should you trust a clustering result?

What will you do differently next Monday?

Validate.
Don't eyeball.



Commit now: before you publish a clustering, run a multi-algorithm + multi-metric + stability check.

Retrieval: name the three validation checks you will run before publishing a clustering.