

Unveiling the Transformative Power of Unsupervised machine learning through Clustering

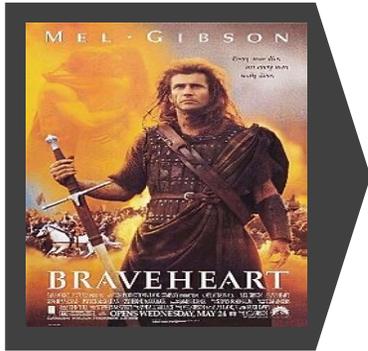
Vishnu S. Pendyala, Ph.D.
San Jose State University

To cite this presentation: Pendyala, V.S. (2025) "Unveiling the Transformative Power of Unsupervised machine learning through Clustering". IEEE Computer Society, Kitchener-Waterloo Chapter Technical Talk, March 31, 2025

©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

How do Streaming services know that these movies can be grouped together?



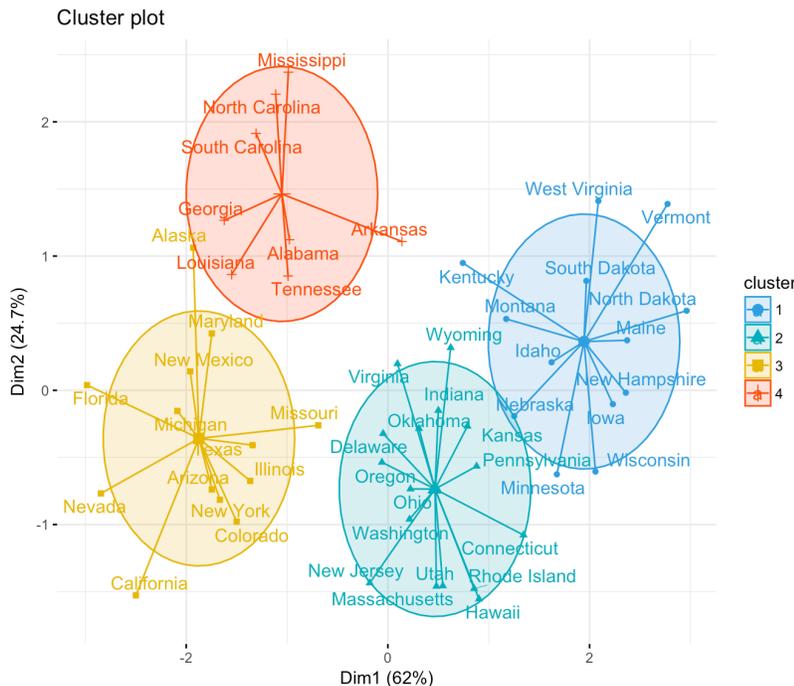


Features are expressed as vectors

Feature	Value
Genre (Action)	8
Genre (Drama)	7
Genre (War)	6
Historical Accuracy	7
Heroism Level	9
Number of Battles	5
Emotional Depth	8
IMDB User Rating	8.4

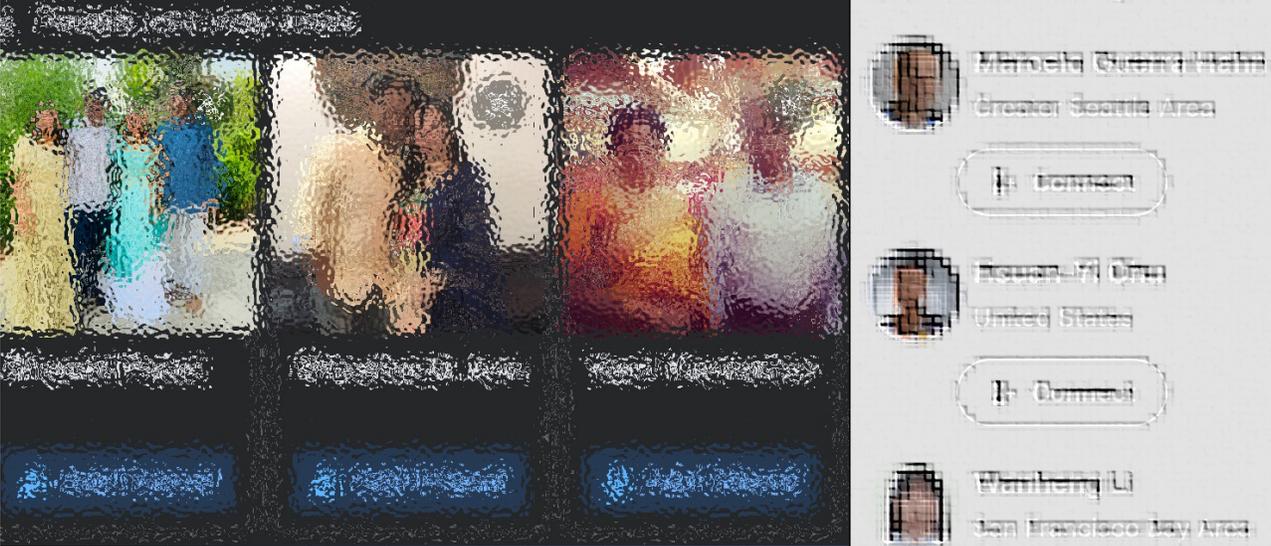
Feature vector $x_i = [8, 7, 6, 7, 9, 5, 8, 8.4]$

©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)



This Photo by Unknown Author is licensed under [CC BY-SA-NC](https://creativecommons.org/licenses/by-sa/4.0/)

- Clustering relies on distances in the feature space
- May not correspond to physical distance
- Distance is a measure for similarity
- Smaller distance => better similarity
- Inter-cluster distances must be maximized
- Intra-cluster distances must be minimized
- No labels (unsupervised)
- Labels => classification



Some applications of Clustering: Social Media

©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

Many applications of clustering

Google Scholar dbscan clustering

Articles About 64,000 results (0.06 sec)

Any time

Since 2025

Since 2024

Since 2021

Custom range...

Sort by relevance

Sort by date

Any type

DBSCAN clustering algorithm based on density

D Deng - 2020 7th international forum on electrical ..., 2020 - ieeexp

... In order to experiment the effect of **DBSCAN** algorithm, this page

DBSCAN algorithm **clustering** on three data sets. These three data

☆ Save Cite Cited by 245 **Related articles** All 2 versions

DBSCAN: Past, present and future

K Khan, SU Rehman, K Aziz, S Fong... - The fifth international ..., 20

... the **DBSCAN** for the purpose of effective **clustering** ... **clustering**

their advantages and limitations. Section IV outlines the critical review

☆ Save Cite Cited by 795 **Related articles** All 6 versions

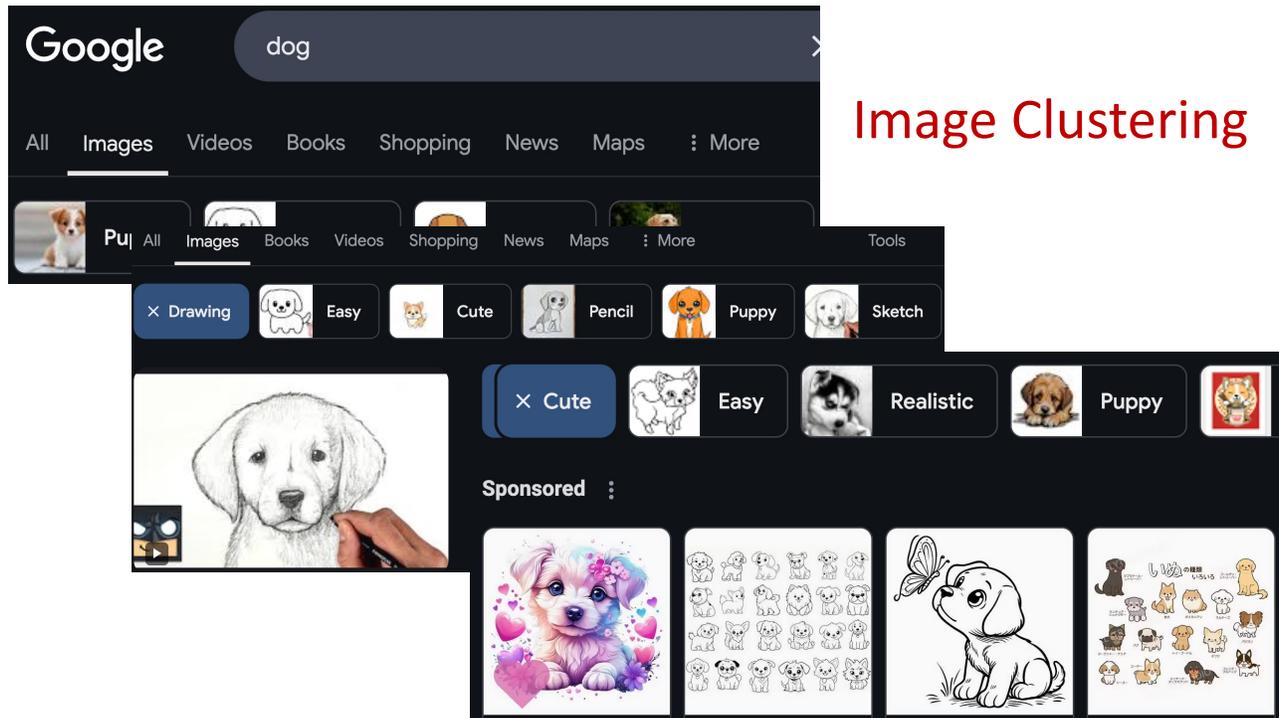


Image Clustering

Classes
vs
Clusters

Classes are defined
and data labeled
manually
(Supervised)

Clusters are deduced
automatically, no
labels
(Unsupervised)

©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

Hyperparameters:
What can you
choose for
clustering data?

Similarity metric: Cosine, Euclidean, Manhattan, Geodesic, ...

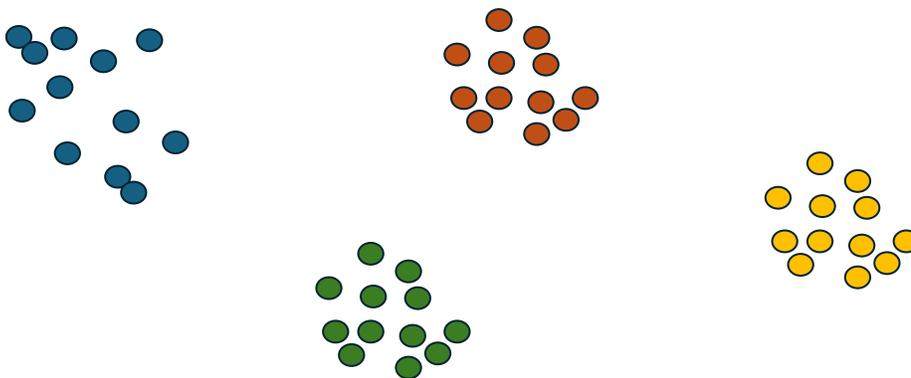
Type of clustering: partitioning (non-overlapping subsets), Hierarchical Clustering (tree-like), Density-Based, Fuzzy Clustering (points to belong to multiple clusters with varying degrees of membership), ...

Clustering algorithm: K-Means, Agglomerative Clustering, Divisive Clustering, DBSCAN, ...

Number of clusters – some algorithms like K-Means need this

©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

How can we detect and form these clusters of data points programmatically?

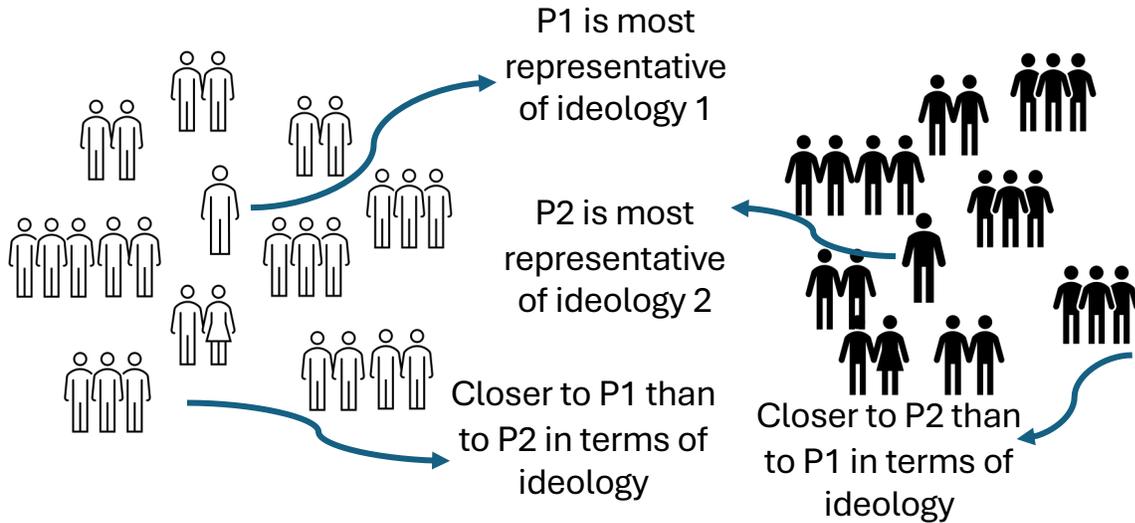


©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

How are clusters formed in real life?

K-means

Political Parties



©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

How are clusters formed in real life?

Neighborhoods

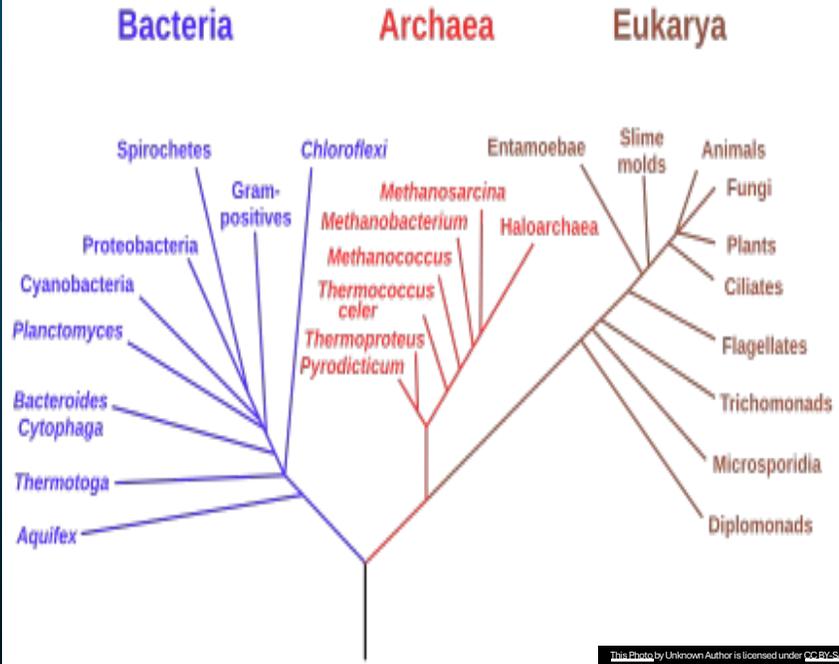


©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

Hierarchical

How are clusters formed in real life?

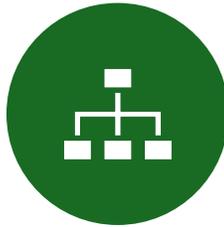
Biological Species



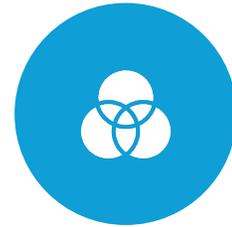
Clustering algorithms for today



K-MEANS



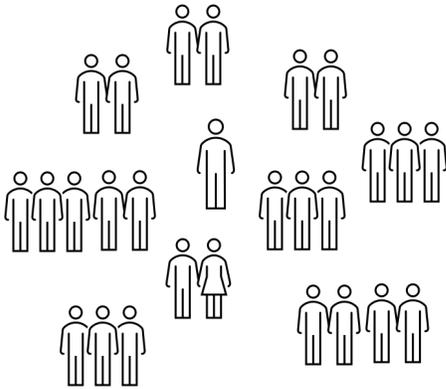
HIERARCHICAL CLUSTERING



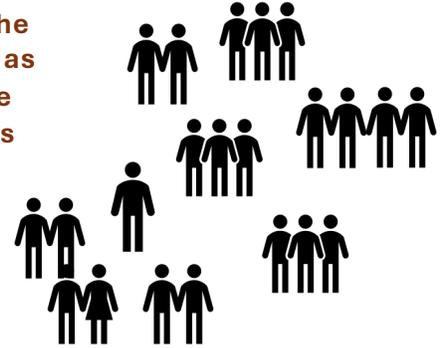
DENSITY-BASED SPATIAL CLUSTERING OF APPLICATIONS WITH NOISE (DBSCAN)

©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

The Goal of Clustering

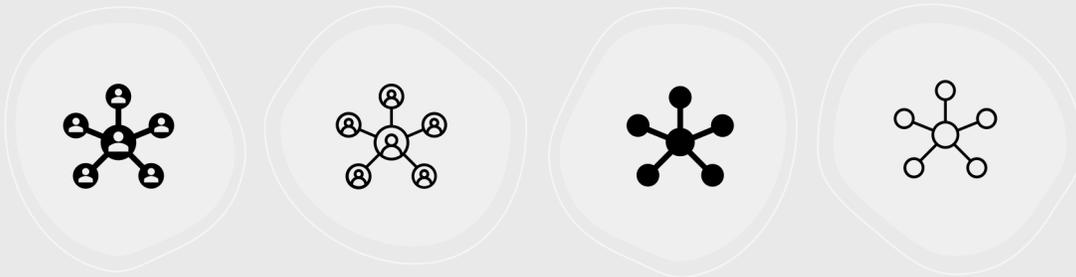


Entities across the clusters must be as **dissimilar** in the feature space as possible



Entities within the clusters must be as **similar** in the feature space as possible

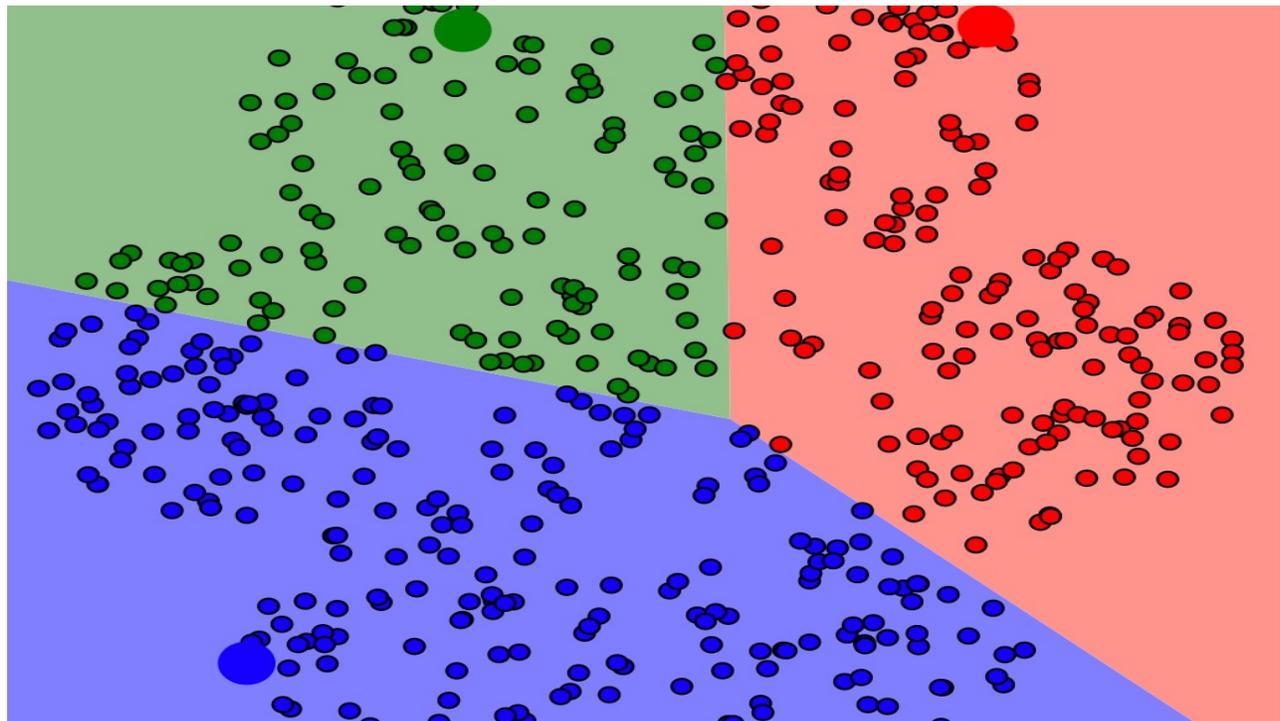
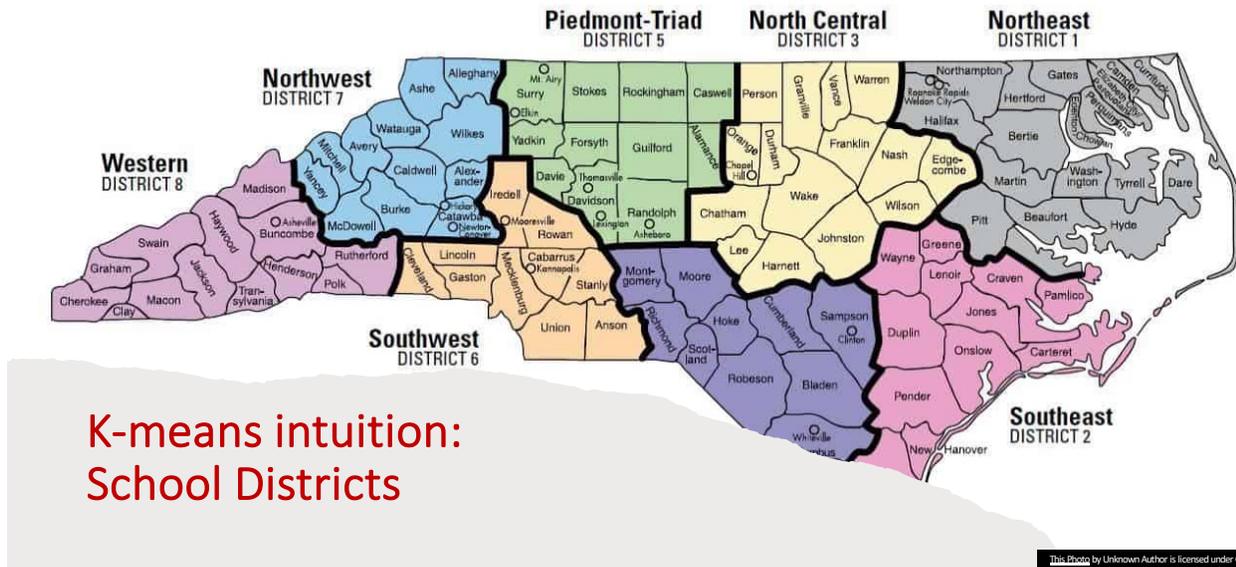
©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)



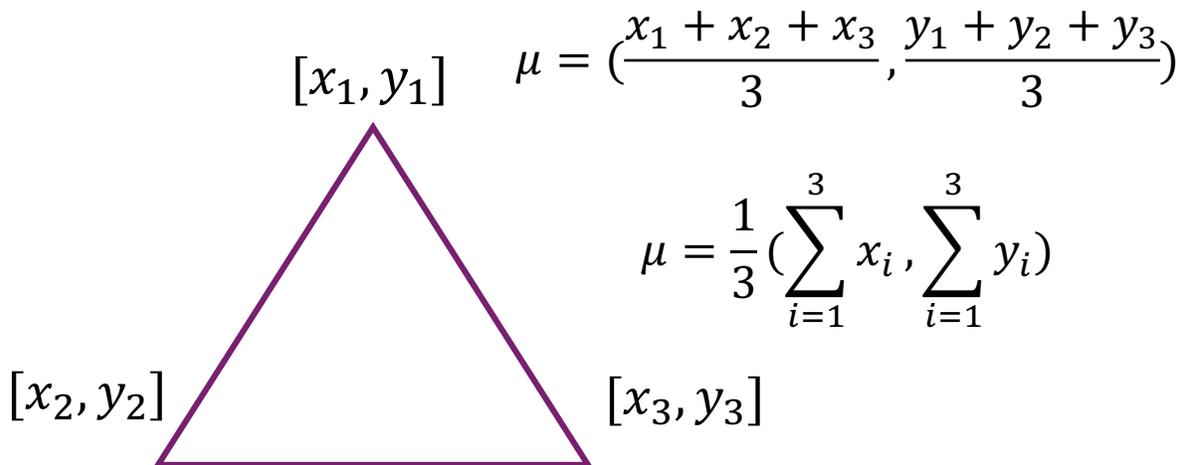
K-Means Clustering

©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

North Carolina State Board of Education Districts



Central to K-means: Centroid

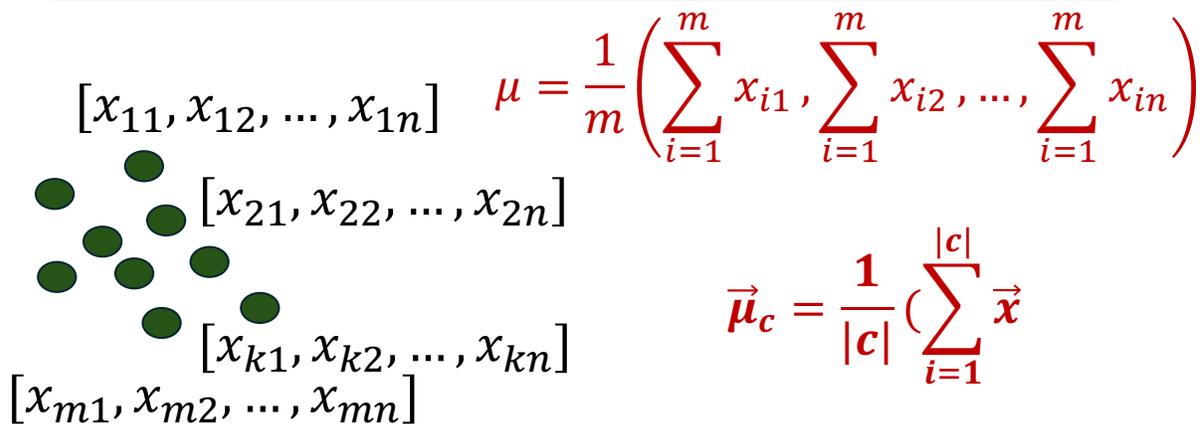


$$\mu = \left(\frac{x_1 + x_2 + x_3}{3}, \frac{y_1 + y_2 + y_3}{3} \right)$$

$$\mu = \frac{1}{3} \left(\sum_{i=1}^3 x_i, \sum_{i=1}^3 y_i \right)$$

©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

Centroid of a cluster



$$\mu = \frac{1}{m} \left(\sum_{i=1}^m x_{i1}, \sum_{i=1}^m x_{i2}, \dots, \sum_{i=1}^m x_{in} \right)$$

$$\vec{\mu}_c = \frac{1}{|c|} \left(\sum_{i=1}^{|c|} \vec{x} \right)$$

m data points

©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

Clustering Hyperparameter: Similarity Metric

Metric	Formula	Properties	Best Used For	Limitations
Euclidean Distance	$\sqrt{\sum_{i=1}^n (x_i - y_i)^2}$	<ul style="list-style-type: none"> - Intuitive - Preserves original space 	<ul style="list-style-type: none"> - Continuous data - Low-dimensional spaces 	<ul style="list-style-type: none"> - Sensitive to outliers - Struggles with high dimensions
Manhattan Distance	$\sum_{i=1}^n x_i - y_i $	<ul style="list-style-type: none"> - Less sensitive to outliers - Computationally efficient 	<ul style="list-style-type: none"> - Grid-like path problems 	<ul style="list-style-type: none"> - Less intuitive - May not capture diagonal relationships
Minkowski Distance	$\left(\sum_{i=1}^n x_i - y_i ^p \right)^{1/p}$	<ul style="list-style-type: none"> - Generalizes Euclidean (p=2) and Manhattan (p=1) - Flexible parameter p 	<ul style="list-style-type: none"> - When optimal distance metric is unknown - Tuning to specific datasets 	<ul style="list-style-type: none"> - Parameter selection can be challenging - Computationally intensive for non-integer p

©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

Distance metrics (continued)

Metric	Formula	Properties	Best Used For	Limitations
Cosine Similarity	$\frac{\sum_{i=1}^n x_i y_i}{\sqrt{\sum_{i=1}^n x_i^2} \sqrt{\sum_{i=1}^n y_i^2}}$	<ul style="list-style-type: none"> - Measures angle, not magnitude - Bounded between -1 and 1 	<ul style="list-style-type: none"> - Text analysis - Recommender systems - When direction matters more than magnitude 	<ul style="list-style-type: none"> - Not a true metric (triangle inequality) - Undefined for zero vectors
Mahalanobis Distance	$\sqrt{(x - y)^T \Sigma^{-1} (x - y)}$	<ul style="list-style-type: none"> - Accounts for correlations - Scale-invariant 	<ul style="list-style-type: none"> - Correlated features - Outlier detection - Classification tasks 	<ul style="list-style-type: none"> - Requires covariance matrix estimation - Computationally expensive

©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

Distance metrics (continued)

Metric	Formula	Properties	Best Used For	Limitations
Hamming Distance	Count of positions where vectors differ	<ul style="list-style-type: none"> - Simple to compute - Natural for categorical data 	<ul style="list-style-type: none"> - Binary / categorical features - Error detection 	<ul style="list-style-type: none"> - Limited to same-length sequences - No concept of magnitude
Jaccard Distance	$1 - \frac{ A \cap B }{ A \cup B }$	<ul style="list-style-type: none"> - Ratio-based - Bounded between 0 and 1 	<ul style="list-style-type: none"> - Set-based problems - Document similarity 	<ul style="list-style-type: none"> - Ignores frequency - Sensitive to small sets
Chebyshev Distance	$\max_i (x_i - y_i)$	<ul style="list-style-type: none"> - Determined by maximum difference - Fast to compute 	<ul style="list-style-type: none"> - Warehouse/path logistics - When worst-case difference matters 	<ul style="list-style-type: none"> - Ignores differences in other dimensions - Sensitive to outliers in single dimension

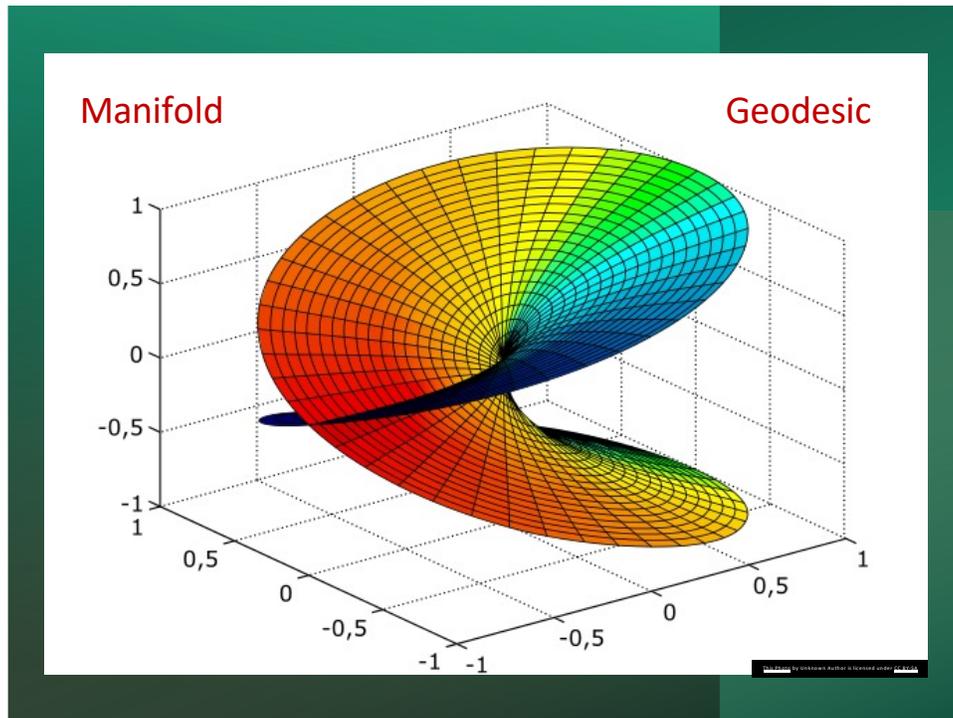
©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](#)

What is the problem with Euclidean distance?

It assumes that the points are all on a hyperplane

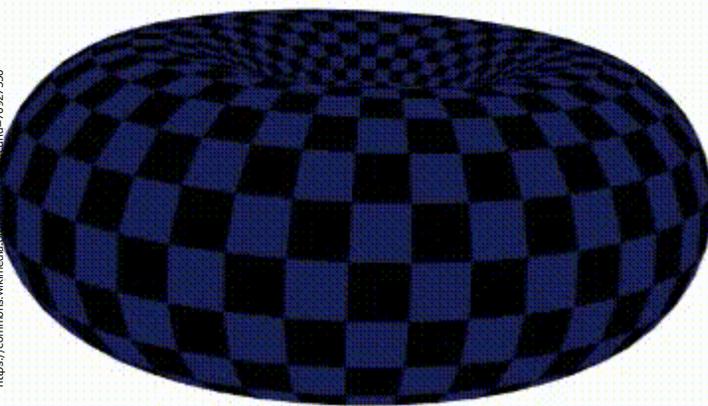


©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](#)



Geodesic –
the
shortest
distance
between
two points
on a
manifold
surface,
honoring
the shape

By Hamishodd11 - Own work, CC BY-SA 4.0
<https://commons.wikimedia.org/wiki/File:Geodesic70927350>



The K-Means Algorithm

- Most popular unsupervised machine learning for partitioning data into K disjoint clusters based on features
- Goal: Minimize within-cluster variance (dissimilarity, measured by distance or sum of squares)

$$Z = \sum_{j=1}^K \sum_{n \in C_j} |x_n - \mu_j|^2$$

- x_n is a vector representing the n^{th} data point, μ_j is the centroid of the data points in the cluster C_j , and $|x_n - \mu_j|^2$ is the Euclidean distance between them.
- Widely used in segmentation problems, quantization, and hidden pattern (structure) recognition

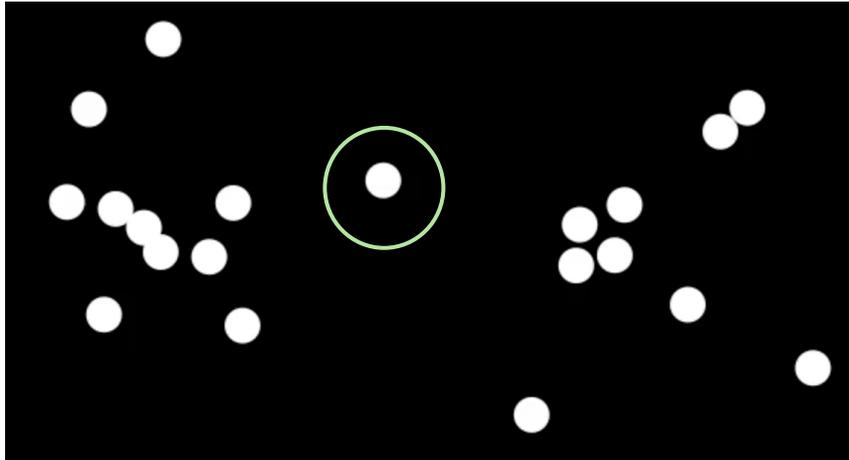
©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

Why is μ_j the centroid?

- Intra-cluster distance $L = \sum (x_n - \mu_j)^2$
- $\frac{\partial L}{\partial \mu} = 2 \sum (x_n - \mu_j) = 0$
- $\Rightarrow \mu_j = \frac{\sum x_n}{P}$, which is the formula for the centroid, where P is the number of points in the cluster
- Proves centroid (or the mean) is the most representative point in the cluster

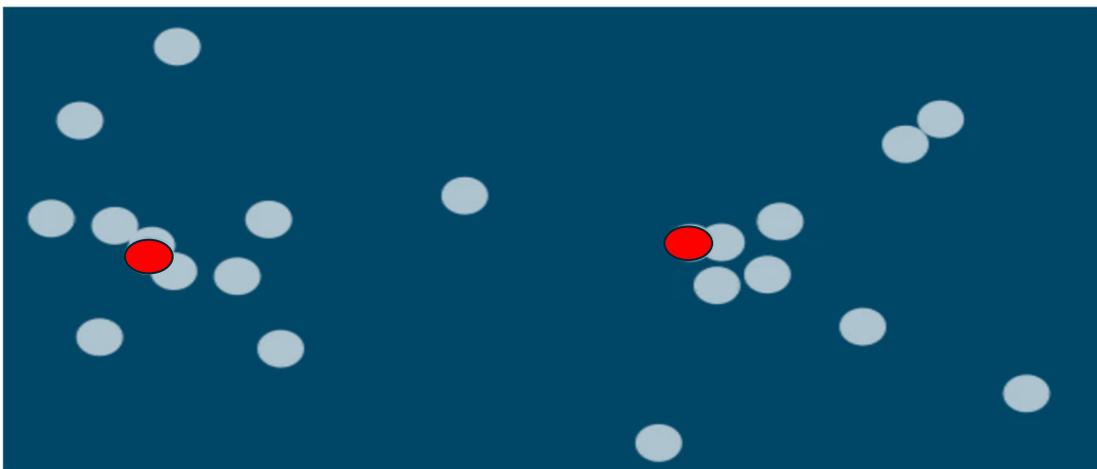
©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

Initial Set of Points



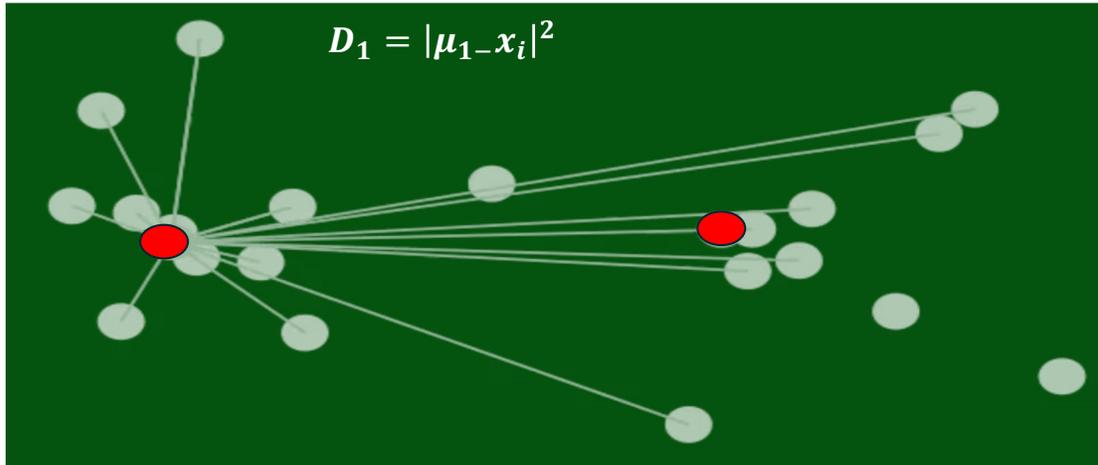
©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

Pseudo-centroids are chosen randomly



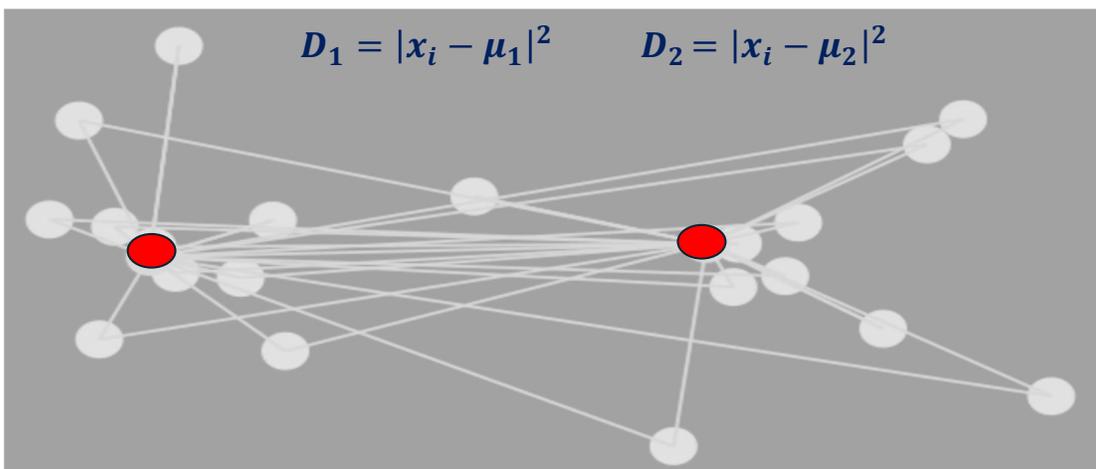
©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

Distances from each data point to the 1st random pseudo-centroid are computed



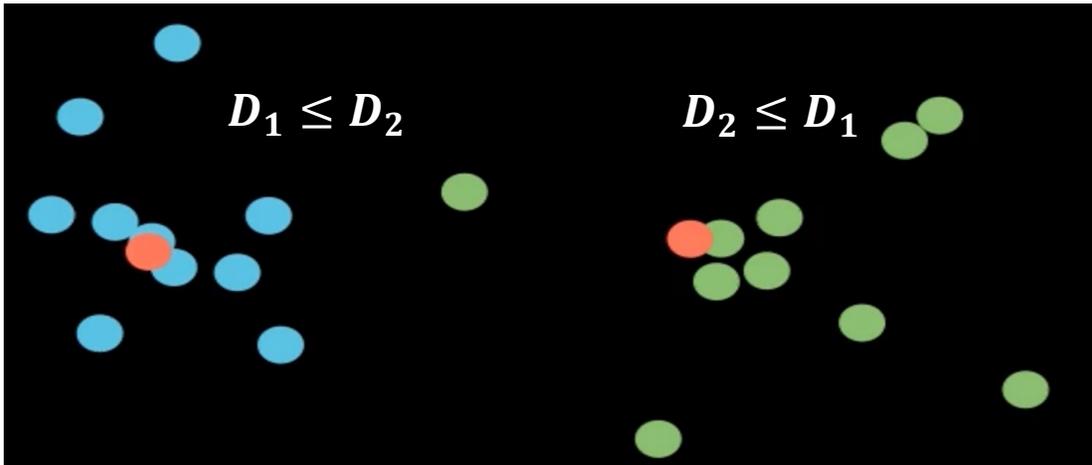
©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

Distances from each data point to the 2nd random pseudo-centroid are computed and compared with the distances from the 1st random centroid



©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

The point in the middle is green since it is closer to the pseudo-centroid representing the green cluster



©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

Recompute the (real) centroids of the clusters

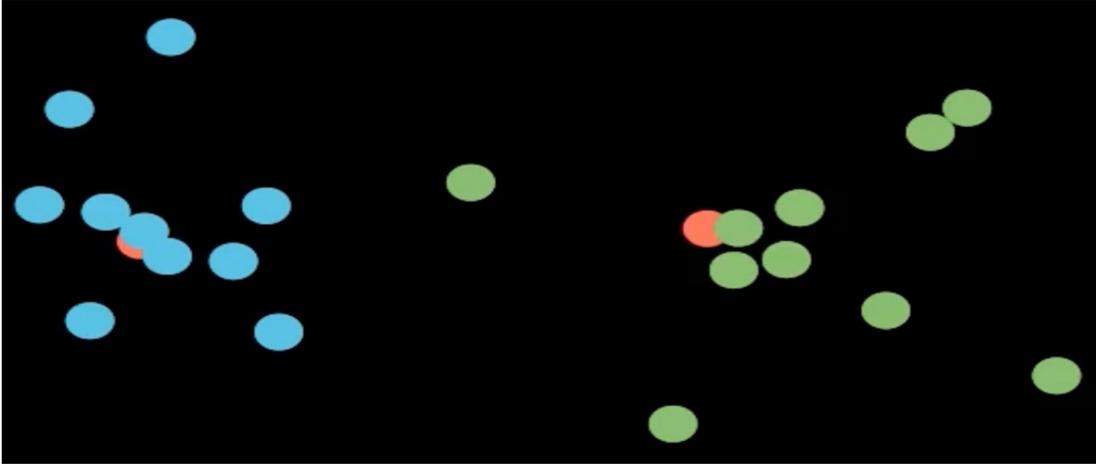
$$\begin{array}{l}
 [x_{11}, x_{12}, \dots, x_{1n}] \\
 [x_{21}, x_{22}, \dots, x_{2n}] \\
 [x_{k1}, x_{k2}, \dots, x_{kn}] \\
 [x_{m1}, x_{m2}, \dots, x_{mn}]
 \end{array}
 \quad
 \mu = \frac{1}{m} \left(\sum_{i=1}^m x_{i1}, \sum_{i=1}^m x_{i2}, \dots, \sum_{i=1}^m x_{in} \right)$$

$$\vec{\mu}_c = \frac{1}{|c|} \left(\sum_{i=1}^{|c|} \vec{x} \right)$$

m data points

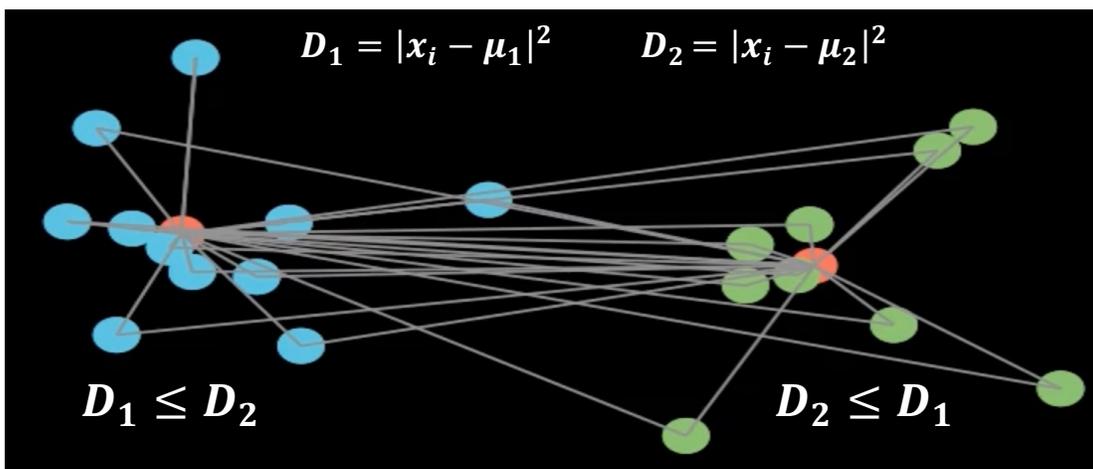
©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

The pseudo-centroids move inside the clusters toward their real centers to become real centroids



©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

Distances are again computed from both (real) centroids and compared => point in the middle now changes to blue



©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

What happens to the objective function with this change to the point in the middle?

$$Z = \sum_{j=1}^K \sum_{n \in C_j} |x_n - \mu_j|^2$$

It reduces!

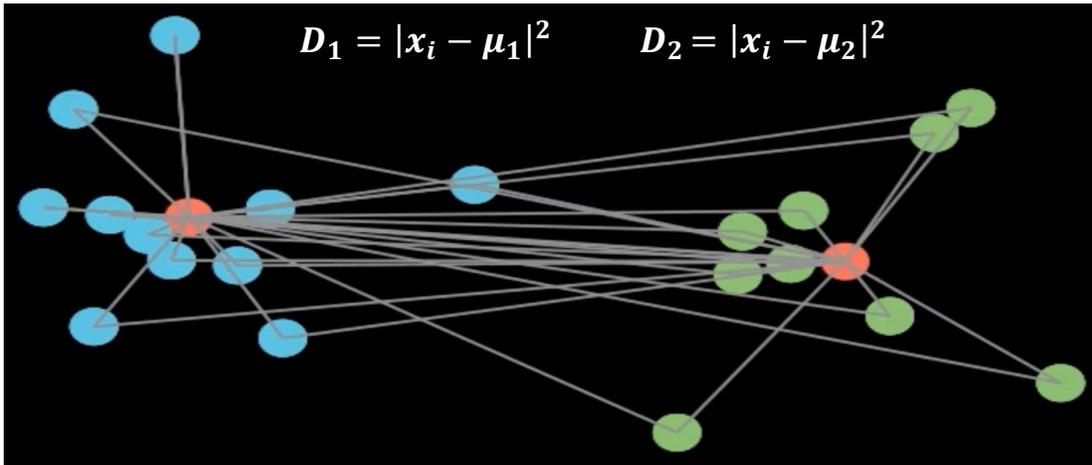
©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

New centroids move further inside



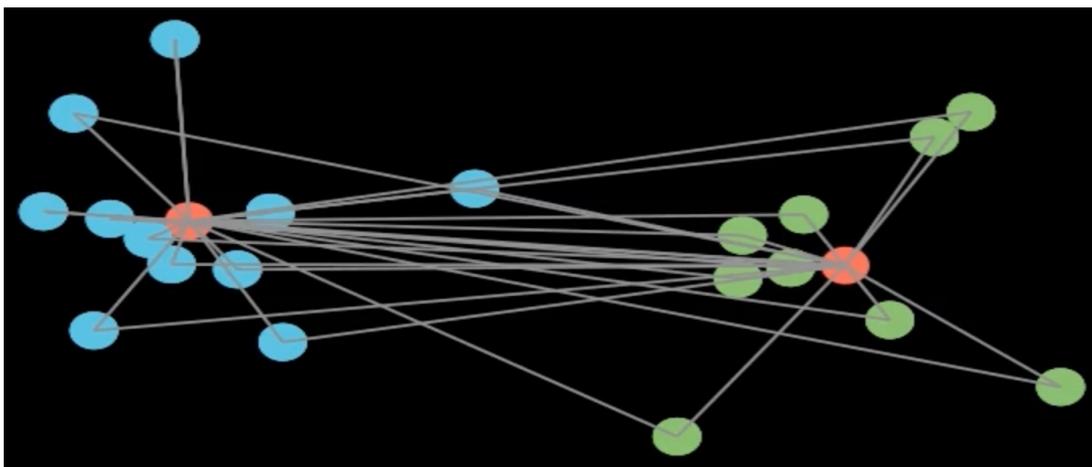
©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

Distances from the new centroids are computed



©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

Based on the distances, the points do not change clusters => convergence

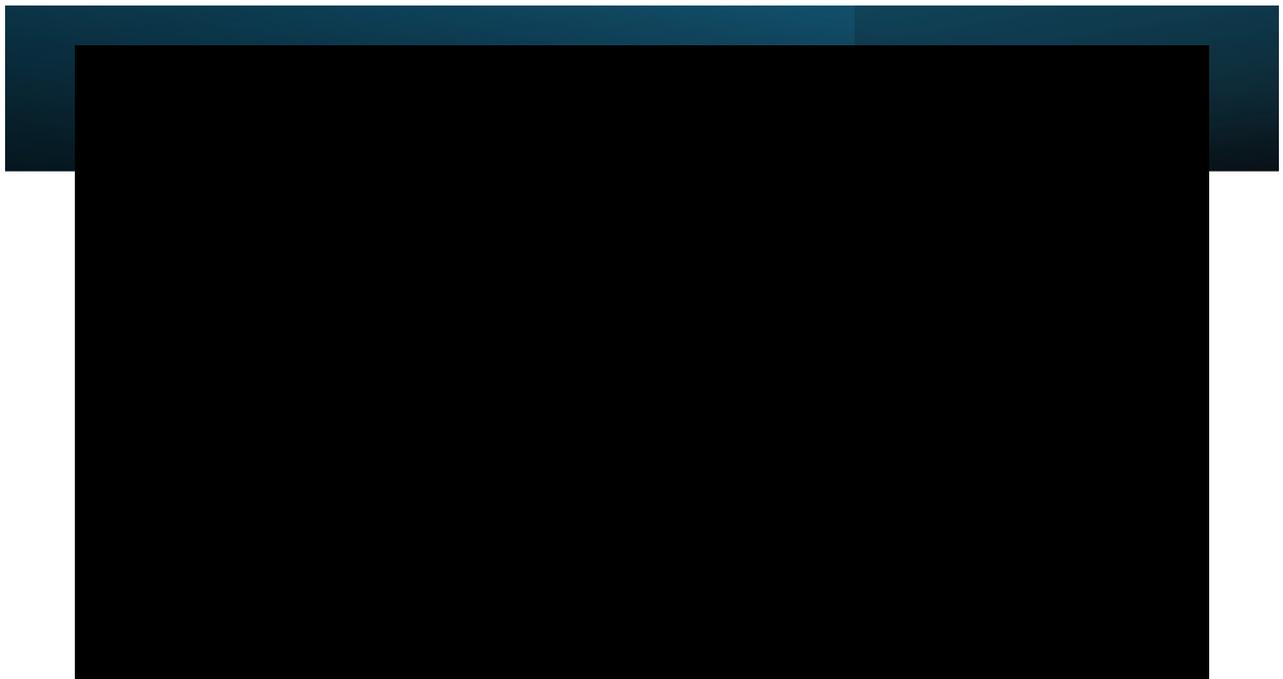


©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

Final clusters after convergence



©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)



©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

K-Means Algorithm

Initialization

- Randomly select K initial centroids: $\mu_1, \mu_2, \dots, \mu_K$
- Initial selection critically impacts final clustering

Repeat

1. Assignment Step

For each data point x :

Assign to closest centroid:

$$\operatorname{argmin}_j \|x - \mu_j\|^2$$

2. Update Step

Recalculate centroids:

$$\mu_i = \left(\frac{1}{|C_i|} \right) \sum_{x \in C_i} x$$

Until

Stopping Conditions:

Centroids no longer move

significantly $\|\mu_{t+1} - \mu_t\| < \epsilon$

(or) No change in cluster assignments

©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

Illustration; 3 Features, k=2; Iteration 1, Step 1: Assignment

x1	x2	x3
1	2	3
1.5	1.8	2.5
5	8	9
8	8	7
1	0.6	1
9	11	12

Random Centroids: Centroid 1: (0, 0, 0) Centroid 2: (10, 10, 10)

x1	x2	x3	Cluster	$D_1 = x_i - \mu_1 ^2$	$D_2 = x_i - \mu_2 ^2$
1	2	3	1	3.742	13.928388
1.5	1.8	2.5	1	3.426368	13.990711
5	8	9	2	13.038405	5.477226
8	8	7	2	13.304135	4.123106
1	0.6	1	1	1.536229	15.822768
9	11	12	2	18.601075	2.449490

$$D_1 = \sqrt{(0-1)^2 + (0-2)^2 + (0-3)^2} = \sqrt{14} = 3.742$$

$$D_2 = \sqrt{(10-1)^2 + (10-2)^2 + (10-3)^2} = \sqrt{81 + 64 + 49} = \sqrt{194} = 13.928$$

©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

Iteration 1: Update Step: Recalculate centroids.

•**Centroid 1:** Mean of (1, 2, 3), (1.5, 1.8, 2.5), (1, 0.6, 1)

$$= ((1+1.5+1)/3, (2+1.8+0.6)/3, (3+2.5+1)/3) = (1.167, 1.467, 2.167)$$

•**Centroid 2:** Mean of (5, 8, 9), (8, 8, 7), (9, 11, 12) = (7.333, 9, 9.333)

x1	x2	x3	Cluster	$D_1 = x_i - \mu_1 ^2$	$D_2 = x_i - \mu_2 ^2$
1	2	3	1	3.742	13.928388
1.5	1.8	2.5	1	3.426368	13.990711
5	8	9	2	13.038405	5.477226
8	8	7	2	13.304135	4.123106
1	0.6	1	1	1.536229	15.822768
9	11	12	2	18.601075	2.449490

©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

Iteration 2: Assignment Step**From Iteration 1**

•**Centroid 1:** (1.167, 1.467, 2.167)

•**Centroid 2:** (7.333, 9, 9.333)

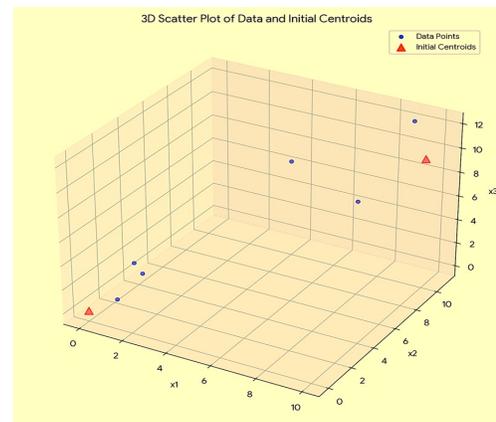
x1	x2	x3	C	$D_1 = x_i - \mu_1 ^2$	$D_2 = x_i - \mu_2 ^2$
1.0	2.0	3.0	1	1.003328	11.367595
1.5	1.8	2.5	1	0.577350	11.513567
5.0	8.0	9.0	2	10.201634	2.560382
8.0	8.0	7.0	2	10.617909	2.624669
1.0	0.6	1.0	1	1.462874	13.420714
9.0	11.0	12.0	2	15.777833	3.726780

From Iteration 2

•**Centroid 1:** (1.167, 1.467, 2.167)

•**Centroid 2:** (7.333, 9, 9.333)

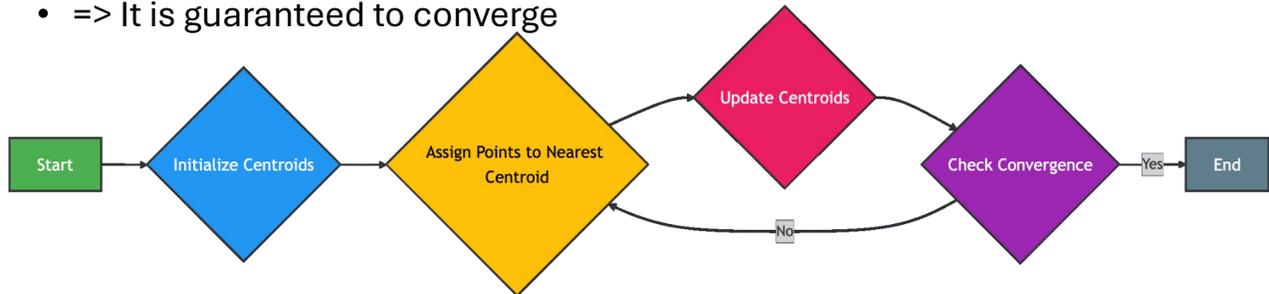
Centroids do not change =>
algorithm converges



©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

K-means Flowchart

- EM is proven to converge
- K-means is a special case of EM
- => It is guaranteed to converge



Expectation

- But convergence does not imply optimality!
- **Bad seeds can result in bad clusters / slow convergence**

©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

How do we know how many clusters?



Representation Learning in lower dimensions and Spectral methods in Machine Learning
... Eigen Decomposition, Eigen Faces, Manifold Learning ...

Dr. Vishnu S. Pendyala, *San Jose State University*

Tuesday, January 16th, 2024, 7:00 pm PT (virtual)
Via Zoom and YouTube Live

Register (Free): <https://r6.ieee.org/scv-cs/representation-learning-in-lower-dimensions-and-spectral-methods-in-machine-learning/>

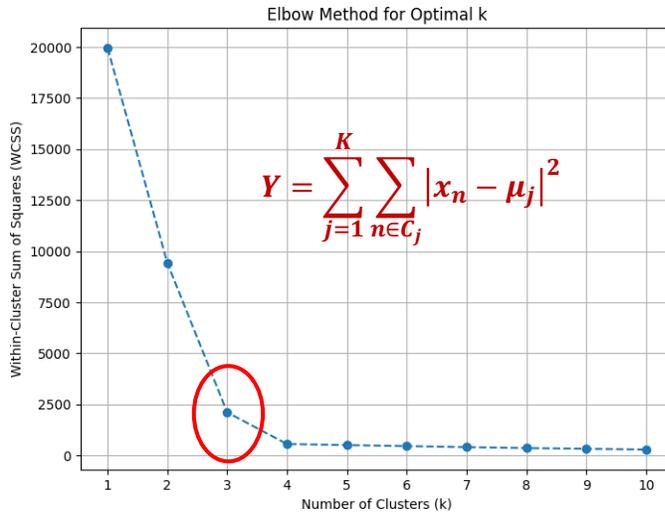
Vishnu S. Pendyala, Chair
 John Delaney, Vice Chair
 Sujata Tibrewala, Secretary
 S.R. Venkatramanan, Treasurer



Santa Clara Valley Chapter

©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

Determining K: Elbow method



K	WCSS
1	19953.76914827803
2	9416.214004352274
3	2110.4125218953286
4	564.9141808210253
5	513.0329042790798
6	462.0960007878173
7	411.3489318727926
8	365.3555992995912
9	332.0364571015475
10	291.5608624106775

©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

Computational complexity of K-Means

Computationally intensive steps:

The algorithm loops for

i: #iterations until convergence

K: #clusters for each cluster

n: #datapoints to compute distance of each point from k centroids

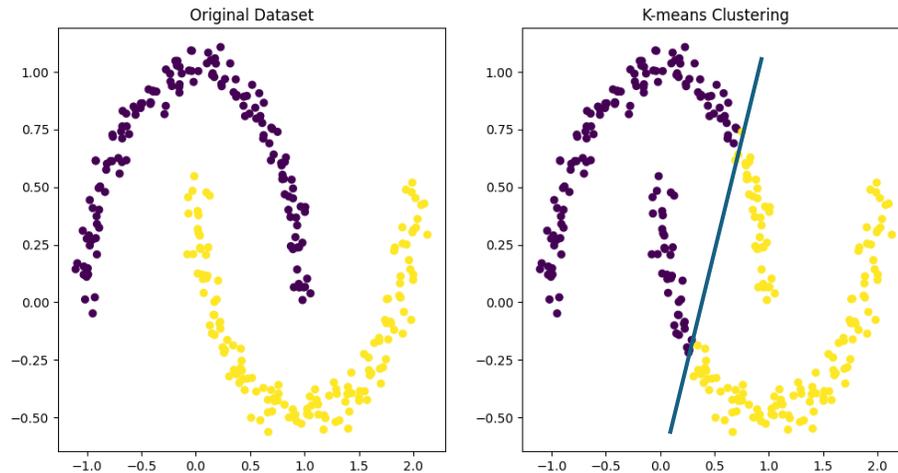
d: #dimensions for computing Euclidean distance from centroid

$$\text{E.g.: } \sqrt{(10 - 1)^2 + (10 - 2)^2 + (10 - 3)^2}$$

Therefore, time complexity: $O(n * K * d * i)$

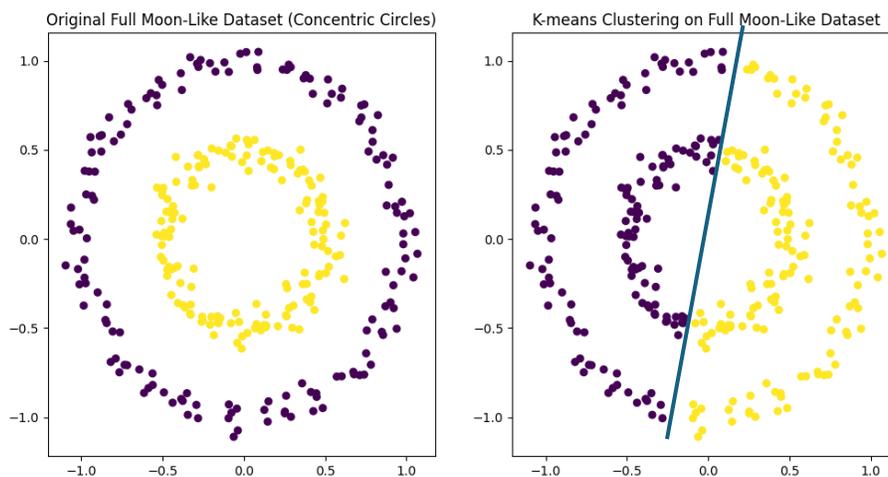
©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

K-Means can only draw linear boundaries!



©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

K-Means can only draw linear boundaries!



©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)



DBSCAN – Key insights

- Clusters data points that are closely packed (density).
- Points that are in low-density regions are flagged as noise.

Simple algorithm based on two hyperparameters to define “dense”

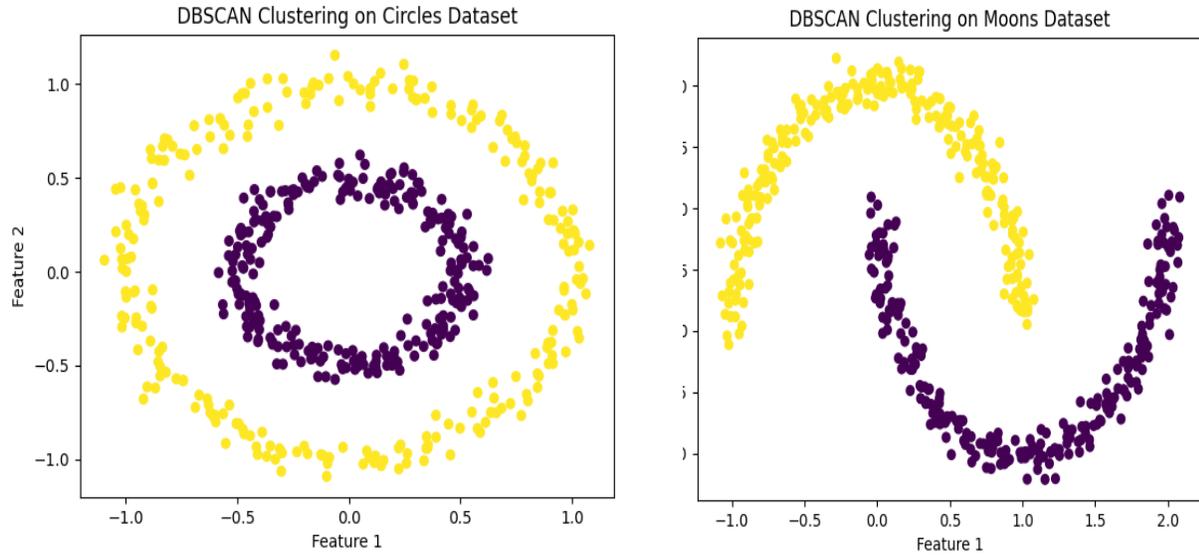
- **Epsilon (ϵ):** Maximum distance between two points for them to be considered as neighbors.
- **MinPts:** Minimum number of points required to form a dense region within a radius ϵ .

Some Applications:

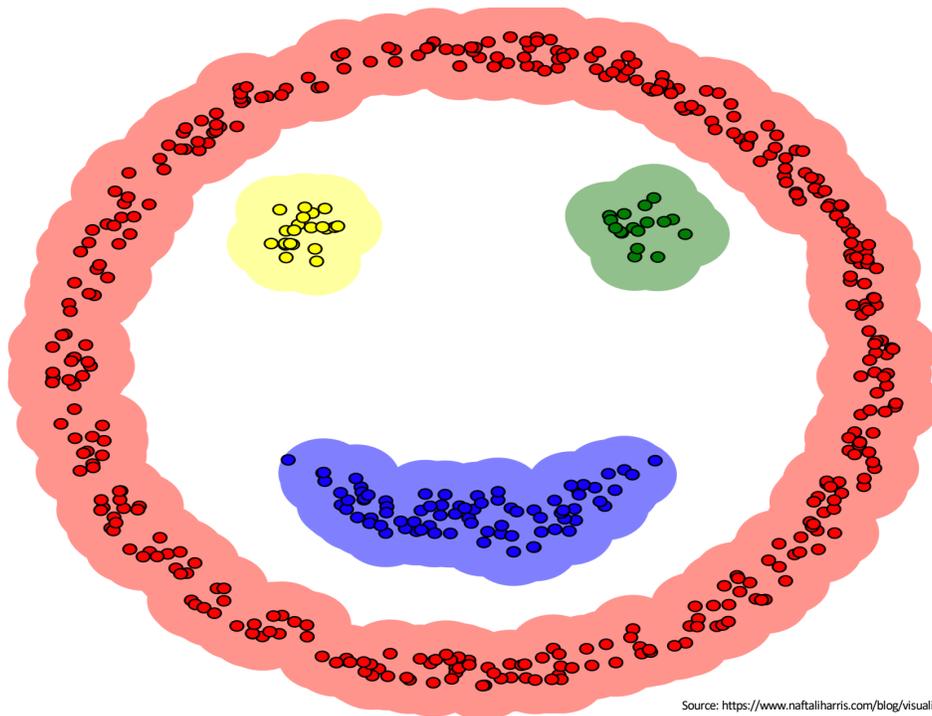
- Geospatial Data: Identifying regions with high population density.
- Anomaly Detection: Detecting fraud or irregular patterns in data.
- Image Segmentation: Identifying regions of interest in images.

©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

DBSCAN can form clusters of varied shapes



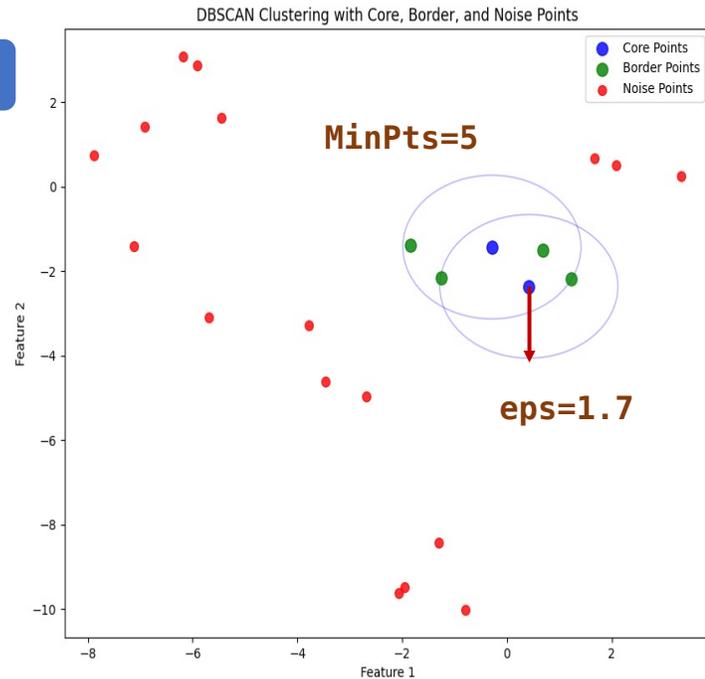
©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)



Source: <https://www.naftaliharris.com/blog/visualizing-dbscan-clustering/>

Core, Border, and Noise Points

- **Core Points:** Points with at least $MinPts$ neighbors within ϵ and will be included in a cluster.
- **Border Points:** Points that are reachable from a core point but do not have enough neighbors to be core points.
- **Noise Points:** Points that are neither core nor border points, often considered outliers.



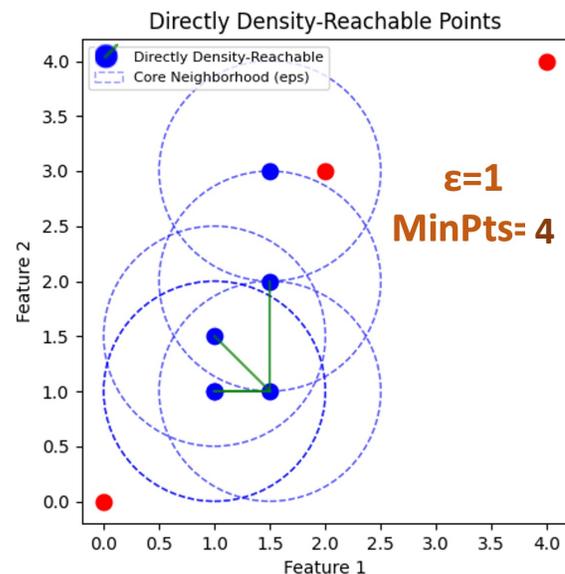
©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

Directly density-reachable

A point p is **directly density-reachable** from point q w.r.t. ϵ and $MinPts$ if

- p is within the ϵ -neighborhood of q : $distance(p, q) \leq \epsilon$ and
- q is a core point:
 $|\{p' \in D \mid distance(q, p') \leq \epsilon\}| \geq MinPts$

Directly density reachability is not symmetric!

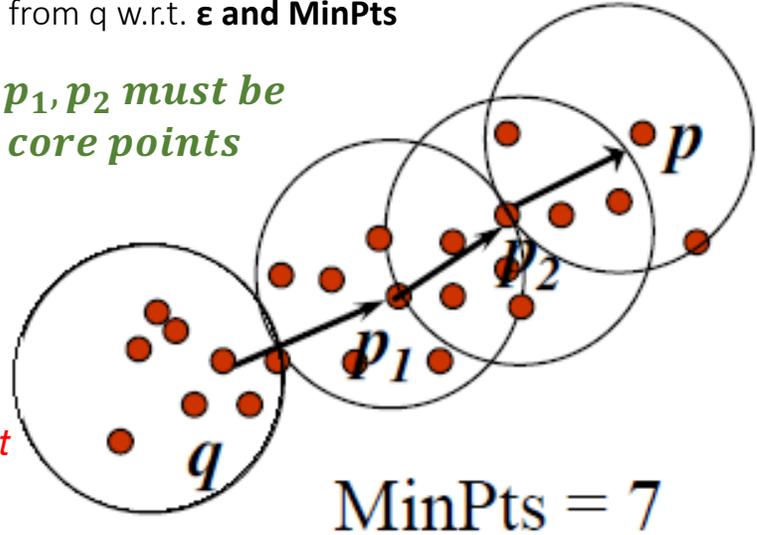


©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

There is a chain of points that are directly density-reachable, starting from q and ending at $p \Rightarrow p$ is density reachable from q w.r.t. ϵ and MinPts

Density Reachable Points

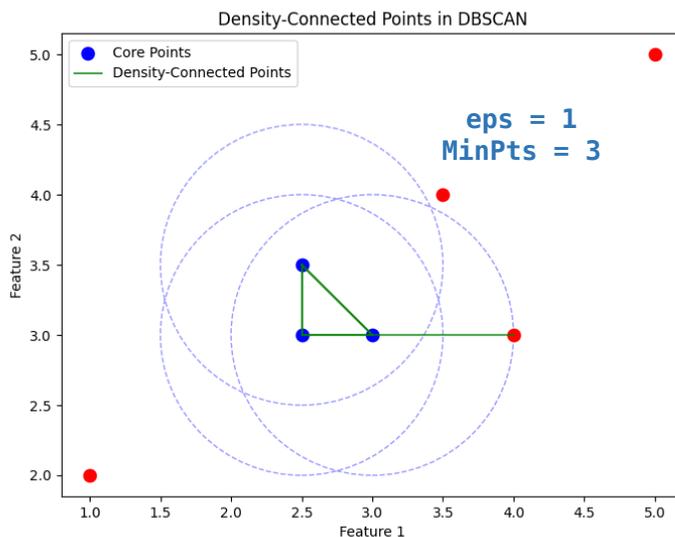
q, p_1, p_2 must be core points



Density reachability is also not symmetric because direct density-reachability is not!

This Photo by Unknown Author is licensed under CC BY-SA

Density Connected



p and q are said to be **density-connected** w.r.t. ϵ and MinPts if there exists a point v such that:

- p is directly density-reachable from v
- q is directly density-reachable from v

Density connected is symmetric!

©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

DBSCAN – the algorithm - initialization

Input: A set of data points, along with two key parameters:

- ϵ (epsilon): The maximum radius of the neighborhood around a point.
- minPts: The minimum number of points required to form a dense region (cluster).

Output: A set of clusters, with some points possibly marked as noise.

For each data point, DBSCAN classifies it into one of three categories:

- Core Point: A point that has at least minPts points (including itself) within its ϵ -neighborhood.
- Border Point: A point that has fewer than minPts points within its ϵ -neighborhood but is within the ϵ -neighborhood of a core point.
- Noise Point: A point that is neither a core point nor a border point.

©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

DBSCAN Algorithm - Forming clusters

If a point is a core point, a new cluster is started, and all points in its ϵ -neighborhood are added to this cluster.

If any of those neighboring points are core points themselves, their ϵ -neighborhoods are recursively processed.

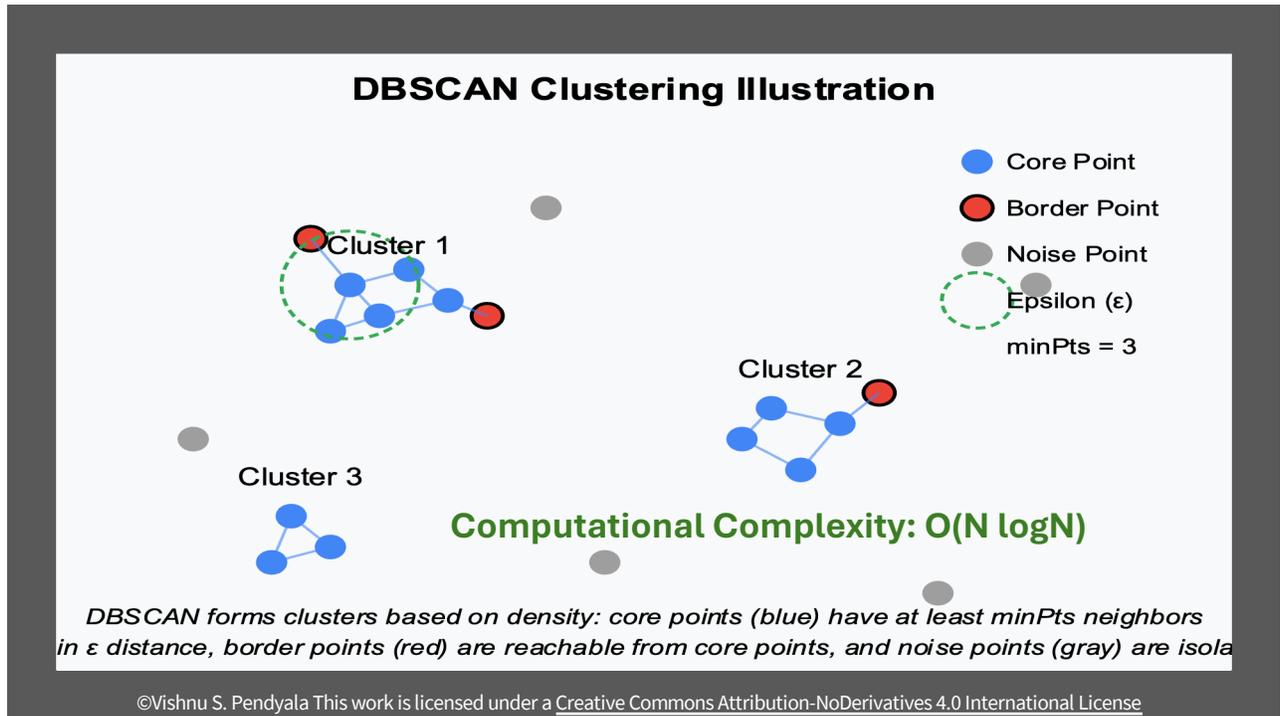
For each core point's ϵ -neighborhood, recursively expand the cluster by checking whether the neighboring points are core points or border points.

Border points are assigned to the cluster of the core point that expanded them.

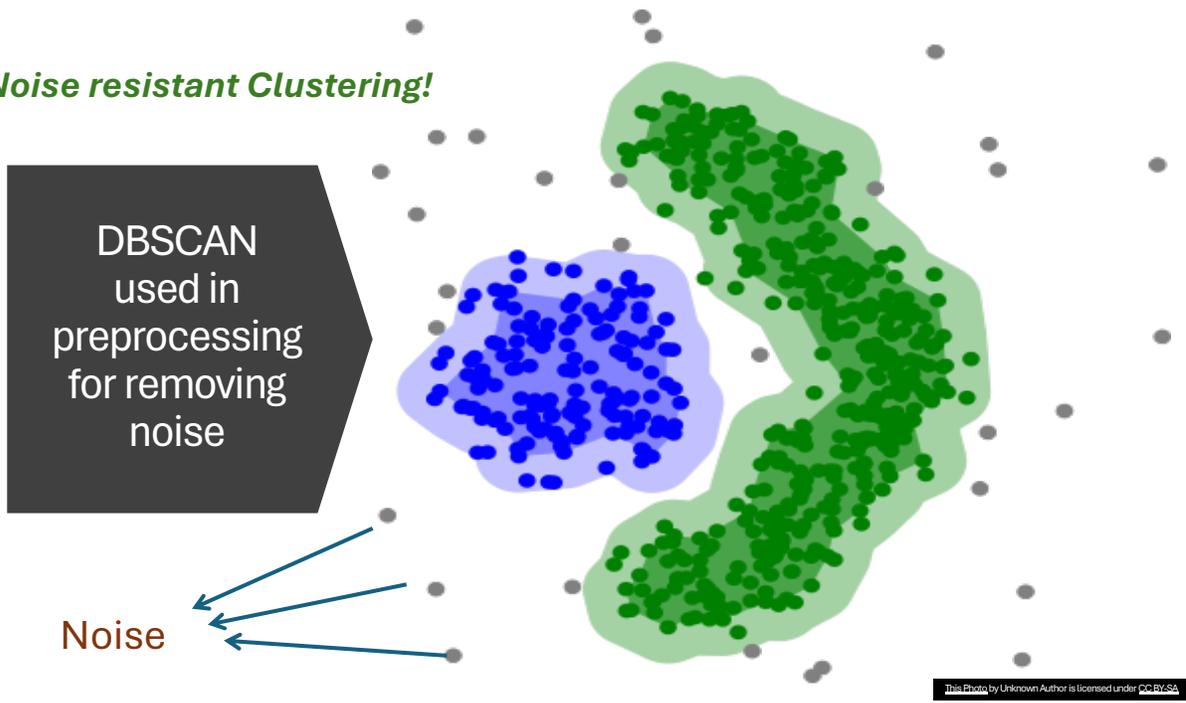
Points that are neither core points nor border points are labeled as noise.

The algorithm stops when all points have been processed (either assigned to a cluster or marked as noise).

©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)



Noise resistant Clustering!



DBSCAN is sensitive to hyperparameters and it is hard to choose the right ones

Figure 8. DBSCAN results for DS1 with MinPts at 4 and Eps at (a) 0.5 and (b) 0.4.

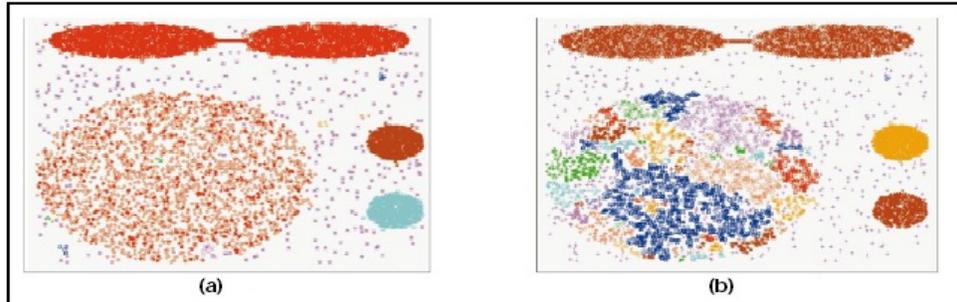
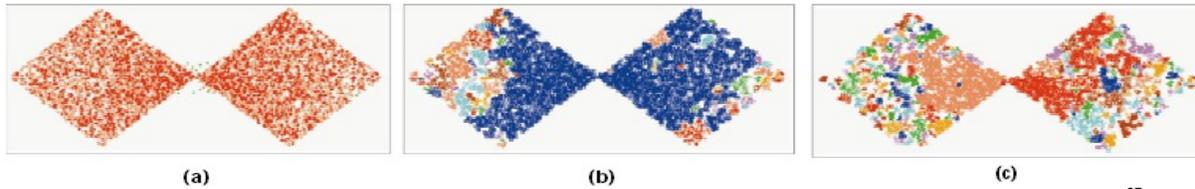


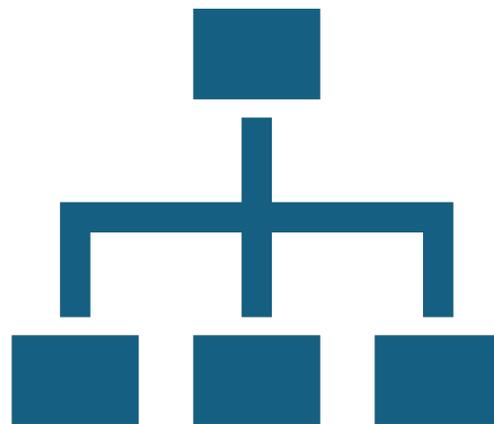
Figure 9. DBSCAN results for DS2 with MinPts at 4 and Eps at (a) 5.0, (b) 3.5, and (c) 3.0.



Source: Data Mining: Concepts and Techniques (3rd ed.)

This Photo by Unknown Author is licensed under CC BY SA

Hierarchical Clustering



Hierarchical Clustering – Key Insights

- Hierarchical clustering organizes data into a hierarchy of nested clusters, visualized as a dendrogram.
- Widely used in Biological sciences, Gene clustering, and Taxonomy creation such as in web catalogs.
- No need to predefine the number of clusters.
- Flexible: Any number of clusters can be obtained by cutting the dendrogram at the desired level.
- Quadratic complexity => Computationally intensive for large datasets

Agglomerative (Bottom-Up):

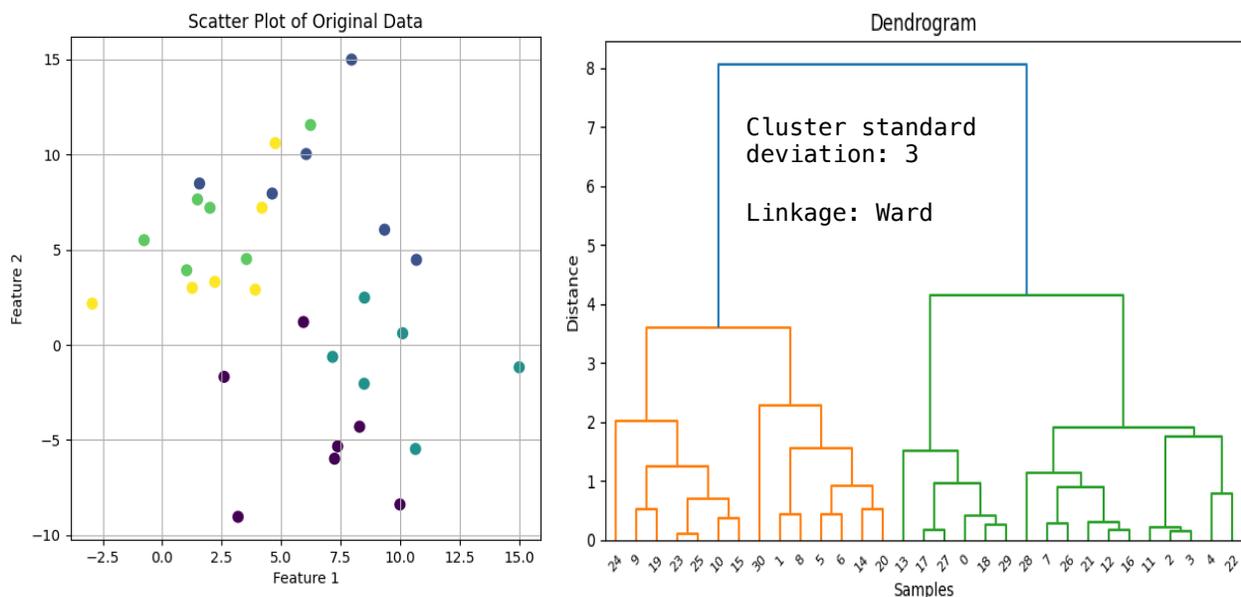
- Start with each point as its own cluster.
- Iteratively merge the closest clusters until one cluster remains.

Divisive (Top-Down):

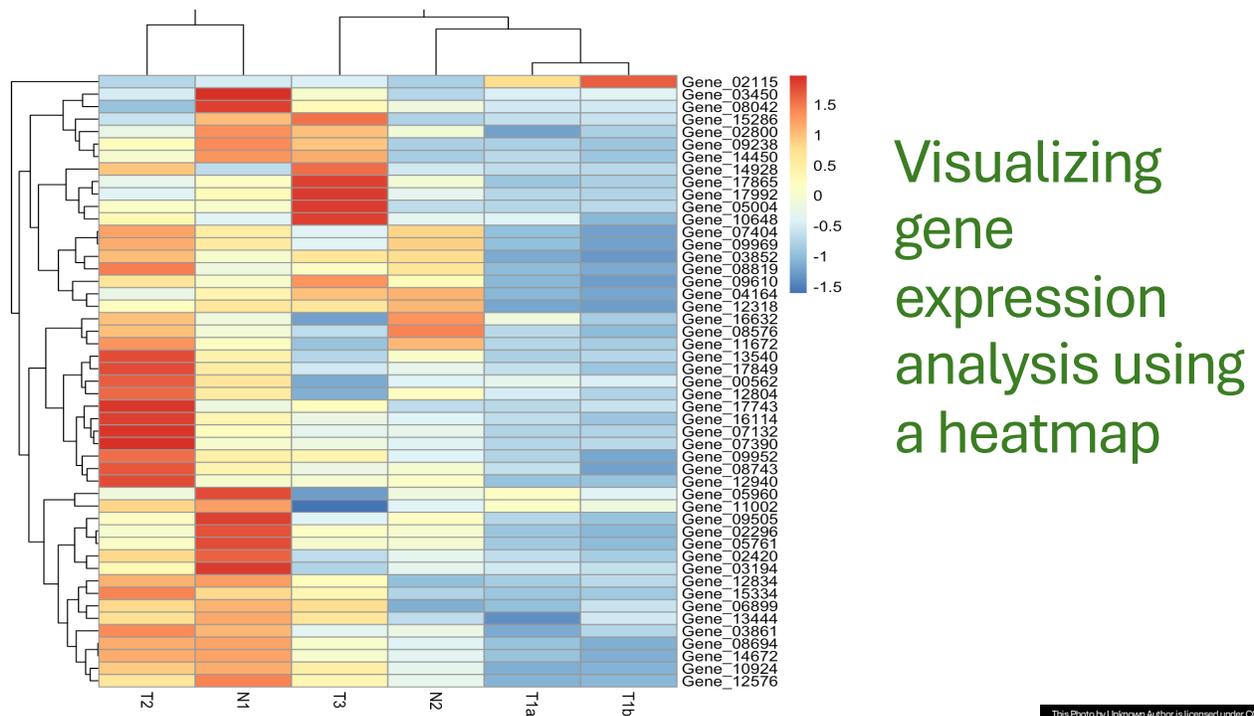
- Start with one all-inclusive cluster.
- Recursively split clusters into smaller groups based on dissimilarity.

©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

Dendrogram for 5 clusters



©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)



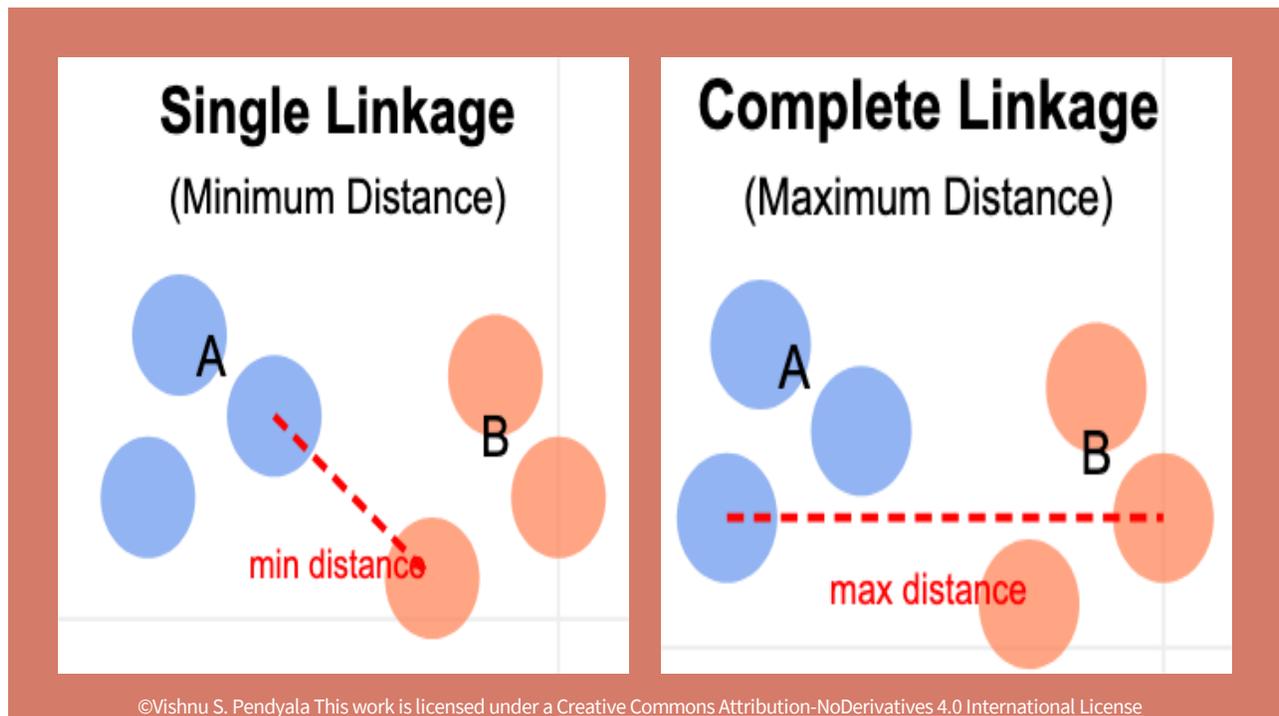
Divisive Clustering Algorithm

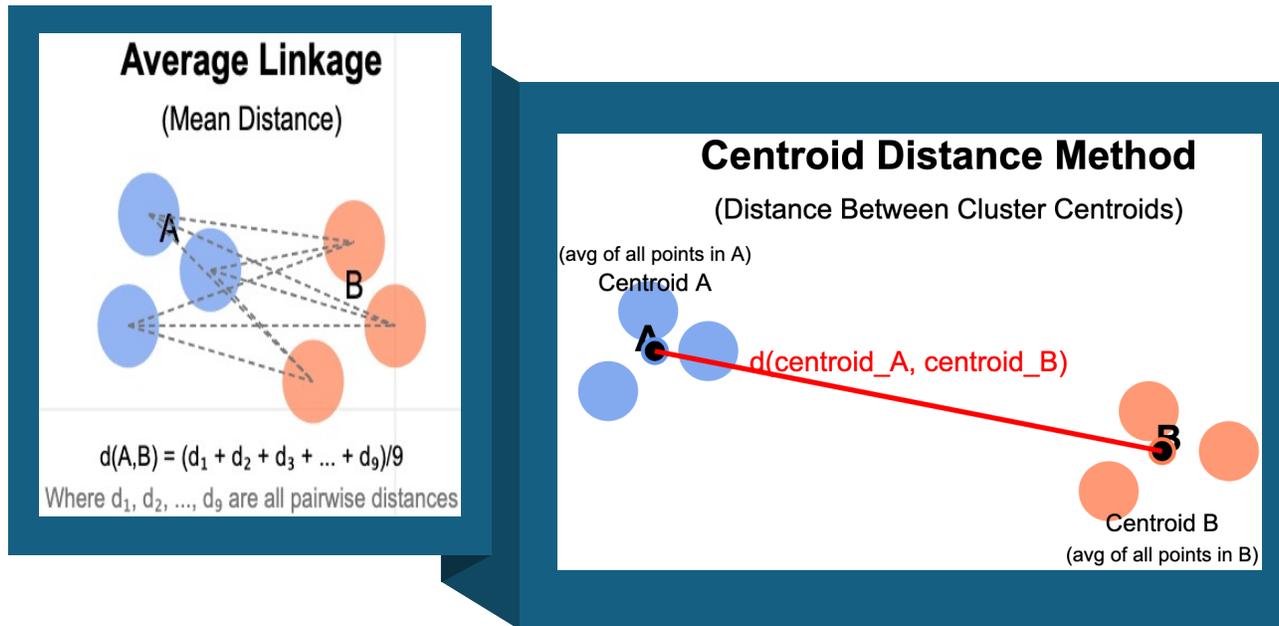
- Start with all points in one cluster.
- Recursively split clusters based on dissimilarity.
- The sequence of splits can be shown using a dendrogram.
- Horizontal cuts to the dendrogram define the number of clusters desired
- Less commonly used due to higher computational demands

Agglomerative Clustering Algorithm

1. Compute the distance matrix between data points.
 2. Let each point be a cluster.
 3. Repeat:
 1. Merge the two closest clusters.
 2. Update the distance matrix.
 4. Stop when only one cluster remains
- **Distance Metrics:**
 - Single linkage (minimum distance).
 - Complete linkage (maximum distance).
 - Average linkage
 - Distance between centroids
 - Ward's method

©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)





©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

Ward's method

Step 1: Original Clusters



©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

Step 2: Calculate Within-Cluster Sum of Squares

For cluster A:

$$\begin{aligned} \text{ESS}(A) &= \sum \|x - \text{centroid}_A\|^2 \\ &= d_1^2 + d_2^2 + d_3^2 \end{aligned}$$

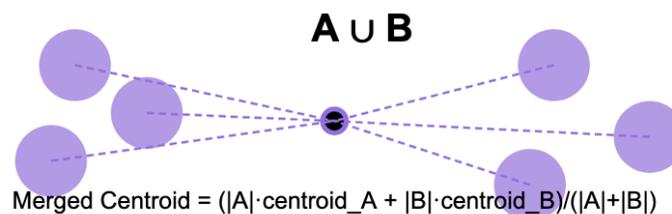
For cluster B:

$$\begin{aligned} \text{ESS}(B) &= \sum \|x - \text{centroid}_B\|^2 \\ &= d_4^2 + d_5^2 + d_6^2 \end{aligned}$$

$$\text{Total Error Sum of Squares} = \text{ESS}(A) + \text{ESS}(B)$$

©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

Step 3: Calculate Error for Hypothetical Merged Cluster



$$\text{Merged Centroid} = (|A| \cdot \text{centroid}_A + |B| \cdot \text{centroid}_B) / (|A| + |B|)$$

$$\begin{aligned} \text{ESS}(A \cup B) &= \sum \|x - \text{centroid}_{\text{merged}}\|^2 \\ &= d_1'^2 + d_2'^2 + d_3'^2 + d_4'^2 + d_5'^2 + d_6'^2 \end{aligned}$$

©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

Step 4: Calculate Ward's Distance

$$\text{Ward's Criterion} = \text{ESS}(A \cup B) - [\text{ESS}(A) + \text{ESS}(B)]$$

$$= \text{Increase in within-cluster sum of squares}$$

©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

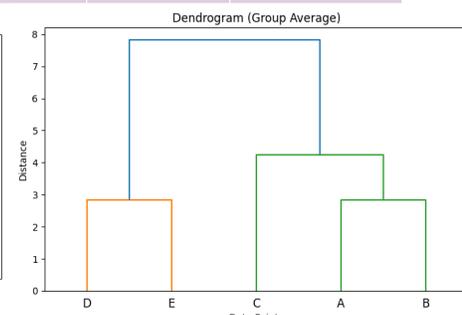
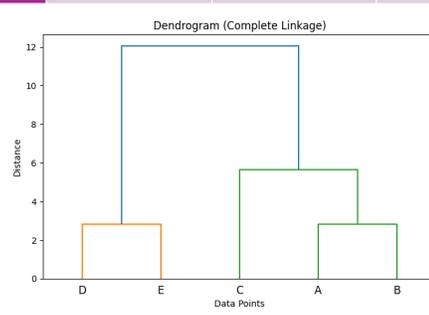
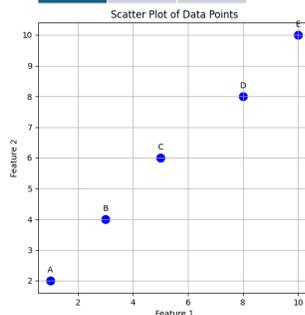
Agglomerative Clustering - Illustration

Dataset

	x1	x2
A	1	2
B	3	4
C	5	6
D	8	8
E	10	10

Proximity Matrix (Euclidean Distance)

	A	B	C	D	E
A	0.000000	2.828427	5.656854	9.219544	12.041595
B	2.828427	0.000000	2.828427	6.403124	9.219544
C	5.656854	2.828427	0.000000	3.605551	6.403124
D	9.219544	6.403124	3.605551	0.000000	2.828427
E	12.041595	9.219544	6.403124	2.828427	0.000000



©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

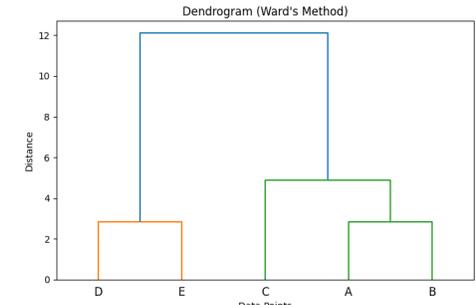
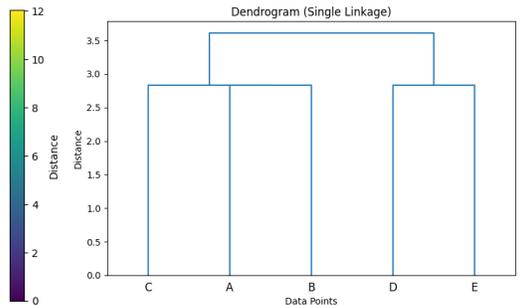
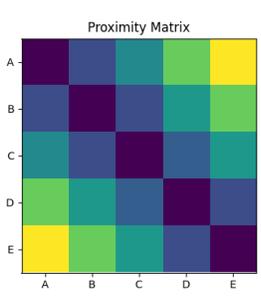
Agglomerative Clustering - Illustration

Dataset

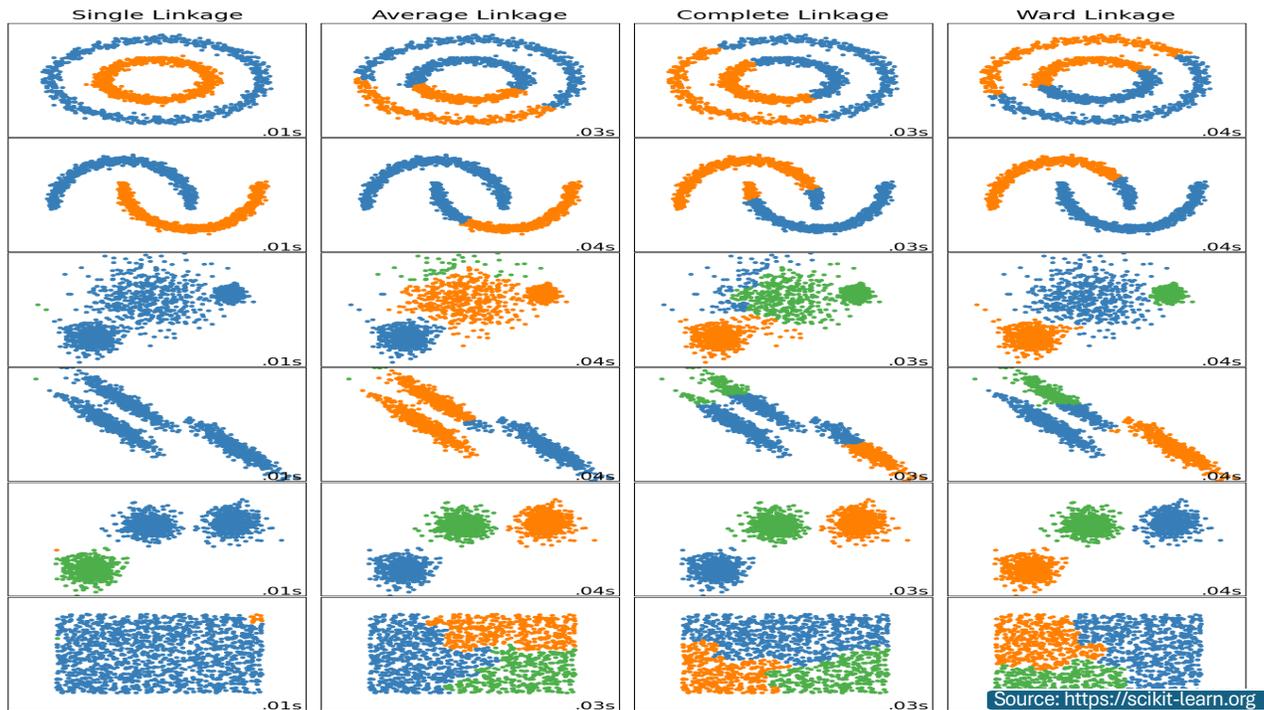
	x1	x2
A	1	2
B	3	4
C	5	6
D	8	8
E	10	10

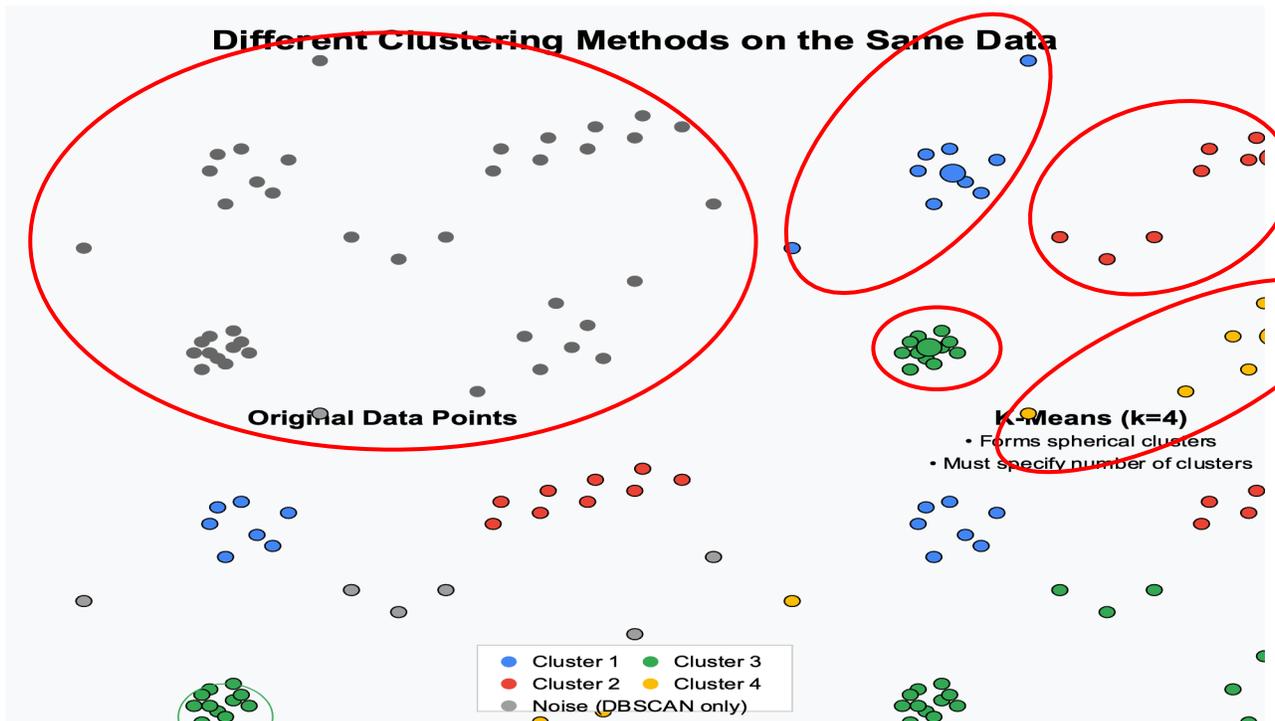
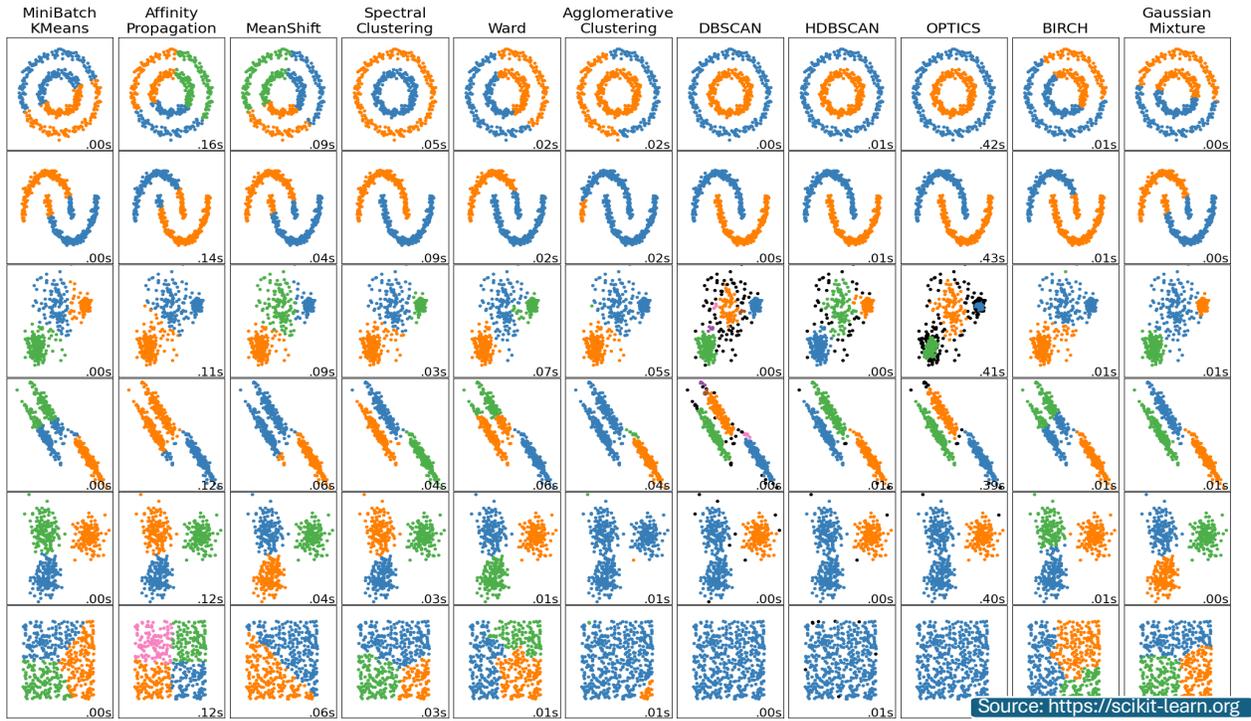
Proximity Matrix (Euclidean Distance)

	A	B	C	D	E
A	0.000000	2.828427	5.656854	9.219544	12.041595
B	2.828427	0.000000	2.828427	6.403124	9.219544
C	5.656854	2.828427	0.000000	3.605551	6.403124
D	9.219544	6.403124	3.605551	0.000000	2.828427
E	12.041595	9.219544	6.403124	2.828427	0.000000



©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)





How do we know which cluster arrangement is the best?

No evaluation metric is perfect; need to depend on heuristics

©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

Evaluating clusters

Metrics: Sum of intra-cluster distances

But what matters is the impact of the clustering on business / application needs!

Final authority in evaluation is the human user

©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](https://creativecommons.org/licenses/by-nd/4.0/)

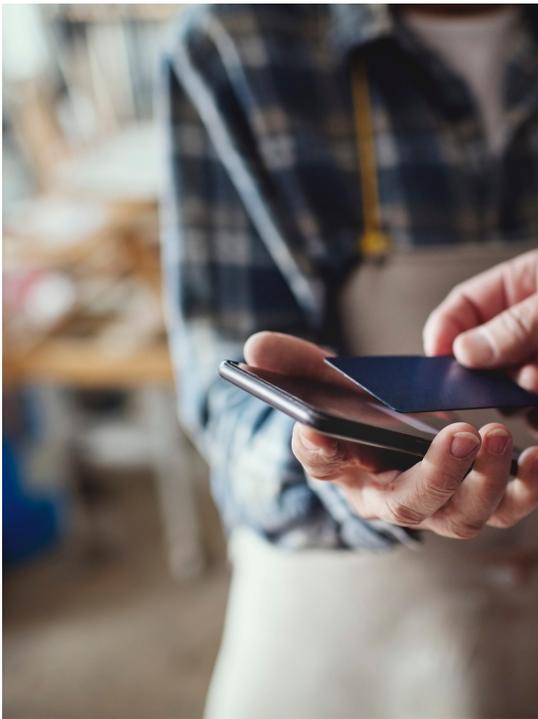


Silhouette Score to evaluate the quality of clustering

Silhouette score 's' = $\frac{b-a}{\max(a,b)}$ where range is [-1, 1] and

- **a:** The average distance between a point and all other points in its own cluster (intra-cluster distance).
- **b:** The average distance between a point and all the points in the nearest cluster it does not belong to (inter-cluster distance).
- Often used for selecting the [optimal number of clusters](#)
- S close to 1: Data point is well-clustered, far from other clusters.
- S close to 0: Data point is on the border of clusters, unclear which cluster it belongs to.
- S close to -1: Data point is probably in the wrong cluster.
- Sensitive to the shape of clusters. May not perform well for non-convex clusters.

©Vishnu S. Pendyala This work is licensed under a [Creative Commons Attribution-NoDerivatives 4.0 International License](#)



Stay in touch!

<https://twitter.com/VishnuPendyala>

<https://www.facebook.com/vishnu.pendyala>

<https://www.instagram.com/vishnupendyala/>

<https://www.threads.net/@vishnupendyala>

<https://www.linkedin.com/in/pendyala/>

Questions

and answers

