

# Lecture 9: Clustering

## Project 7

Να φορτώσετε ένα dataset της επιλογής σας, και να χρησιμοποιήσετε τουλάχιστον 2 Clustering αλγορίθμους (συνιστάται η χρήση της βιβλιοθήκης **scikit** για Python) για την εκπαίδευση ενός μοντέλου μηχανικής μάθησης.

Κατά τη διάρκεια της εφαρμογής του αλγορίθμου, να γίνουν πειράματα τα οποία θα αντιστοιχούν στο συνδυασμό παραμέτρων που βρίσκονται στη στήλη **Parameters** του παρακάτω πίνακα σύγκρισης αλγορίθμων συσταδοποίησης ([source](#)). Να αξιολογηθεί το μοντέλο κάνοντας χρήση της μετρικής **Silhouette** και να παραδοθούν τα αποτελέσματα σε output της επιλογής σας (συμπεριλαμβανομένου των παραμέτρων που χρησιμοποιήθηκαν και των αντίστοιχων τιμών τους).

**Σημείωση:** Κατά την υποβολή της εργασίας, είναι απαραίτητο να συμπεριληφθεί και ο κώδικας που χρησιμοποιήθηκε. Για διευκόλυνση, επισυνάπτεται το αρχείο **Clustering\_Template.py** μέσα στο οποίο μπορεί να συμπληρωθεί ο απαραίτητος κώδικας σε Python.

Method name	Parameters	Scalability	Usecase	Geometry (metric used)
K-Means	number of clusters	Very large <code>n_samples</code> , medium <code>n_clusters</code> with <code>MiniBatch</code> code	General-purpose, even cluster size, flat geometry, not too many clusters	Distances between points
Affinity propagation	damping, sample preference	Not scalable with <code>n_samples</code>	Many clusters, uneven cluster size, non-flat geometry	Graph distance (e.g. nearest-neighbor graph)
Mean-shift	bandwidth	Not scalable with <code>n_samples</code>	Many clusters, uneven cluster size, non-flat geometry	Distances between points
Spectral clustering	number of clusters	Medium <code>n_samples</code> , small <code>n_clusters</code>	Few clusters, even cluster size, non-flat geometry	Graph distance (e.g. nearest-neighbor graph)
Ward hierarchical clustering	number of clusters	Large <code>n_samples</code> and <code>n_clusters</code>	Many clusters, possibly connectivity constraints	Distances between points
Agglomerative clustering	number of clusters, linkage type, distance	Large <code>n_samples</code> and <code>n_clusters</code>	Many clusters, possibly connectivity constraints, non Euclidean distances	Any pairwise distance
DBSCAN	neighborhood size	Very large <code>n_samples</code> , medium <code>n_clusters</code>	Non-flat geometry, uneven cluster sizes	Distances between nearest points
Gaussian mixtures	many	Not scalable	Flat geometry, good for density estimation	Mahalanobis distances to centers
Birch	branching factor, threshold, optional global clusterer.	Large <code>n_clusters</code> and <code>n_samples</code>	Large dataset, outlier removal, data reduction.	Euclidean distance between points