

Clustering Algorithms: Centroid-based

We talked about K-Means Clustering ...

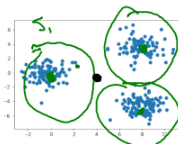
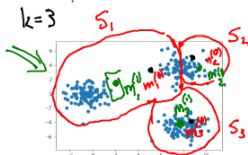
k-Means Clustering Algorithm

- **Initialization:** Choose k and k points randomly in the data to be the initial means (centroids) $m_1^{(0)}, \dots, m_k^{(0)}$.
- **Assignment:** assign each point to the cluster with the nearest mean:

$$\| \hat{x} \|^2 = \left(\sum_{j=1}^k (x_j)^2 \right) \Rightarrow \| \hat{x} - \hat{y} \|^2 = \sum_{j=1}^k (x_j - y_j)^2$$

$$S_i^{(t)} = \{x_p : \|x_p - m_i^{(t)}\|^2 \leq \|x_p - m_j^{(t)}\|^2 \forall j, 1 \leq j \leq k\}$$
- **Update:** update the mean

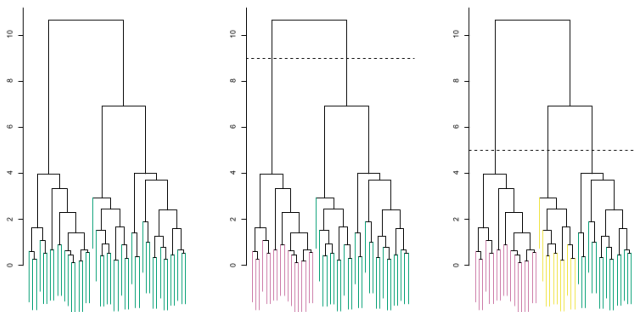
$$m_i^{(t+1)} = \frac{1}{|S_i^{(t)}|} \sum_{x_j \in S_i^{(t)}} x_j$$



There are other approaches ...

Clustering Algorithms: Hierarchical

Hierarchical algorithms proceed by finding nearest neighbors iteratively and then making choices about the number of clusters:



Clustering Algorithms: Hierarchical

The algorithm:

Algorithm 12.3 *Hierarchical Clustering*

1. Begin with n observations and a measure (such as Euclidean distance) of all the $\binom{n}{2} = n(n-1)/2$ pairwise dissimilarities. Treat each observation as its own cluster.
 2. For $i = n, n-1, \dots, 2$:
 - (a) Examine all pairwise inter-cluster dissimilarities among the i clusters and identify the pair of clusters that are least dissimilar (that is, most similar). Fuse these two clusters. The dissimilarity between these two clusters indicates the height in the dendrogram at which the fusion should be placed.
 - (b) Compute the new pairwise inter-cluster dissimilarities among the $i-1$ remaining clusters.
-

The question is: what does it mean to be the nearest neighbor to a “cluster” of points?

Clustering Algorithms: Hierarchical

The Linkages:

| <i>Linkage</i> | <i>Description</i> |
|----------------|---|
| Complete | Maximal intercluster dissimilarity. Compute all pairwise dissimilarities between the observations in cluster A and the observations in cluster B, and record the <i>largest</i> of these dissimilarities. |
| Single | Minimal intercluster dissimilarity. Compute all pairwise dissimilarities between the observations in cluster A and the observations in cluster B, and record the <i>smallest</i> of these dissimilarities. Single linkage can result in extended, trailing clusters in which single observations are fused one-at-a-time. |
| Average | Mean intercluster dissimilarity. Compute all pairwise dissimilarities between the observations in cluster A and the observations in cluster B, and record the <i>average</i> of these dissimilarities. |
| Centroid | Dissimilarity between the centroid for cluster A (a mean vector of length p) and the centroid for cluster B. Centroid linkage can result in undesirable <i>inversions</i> . |

What does SKLearn do?

Clustering Algorithms: Density-Based

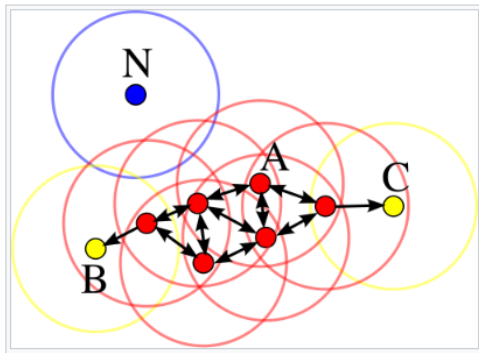
From Wikipedia . . .

Consider a set of points in some space to be clustered. Let ϵ be a parameter specifying the radius of a neighborhood with respect to some point. For the purpose of DBSCAN clustering, the points are classified as *core points*, (*directly-*) *reachable points* and *outliers*, as follows:

- A point p is a *core point* if at least minPts points are within distance ϵ of it (including p).
- A point q is *directly reachable* from p if point q is within distance ϵ from core point p . Points are only said to be directly reachable from core points.
- A point q is *reachable* from p if there is a path p_1, \dots, p_n with $p_1 = p$ and $p_n = q$, where each p_{i+1} is directly reachable from p_i . Note that this implies that the initial point and all points on the path must be core points, with the possible exception of q .
- All points not reachable from any other point are *outliers* or *noise points*.

Now if p is a core point, then it forms a *cluster* together with all points (core or non-core) that are reachable from it. Each cluster contains at least one core point; non-core points can be part of a cluster, but they form its "edge", since they cannot be used to reach more points.

Clustering Algorithms: Density-Based



A “cluster” has two properties:

1. All points within the cluster are mutually density-connected.
2. If a point is density-reachable from some point of the cluster, it is part of the cluster as well.