

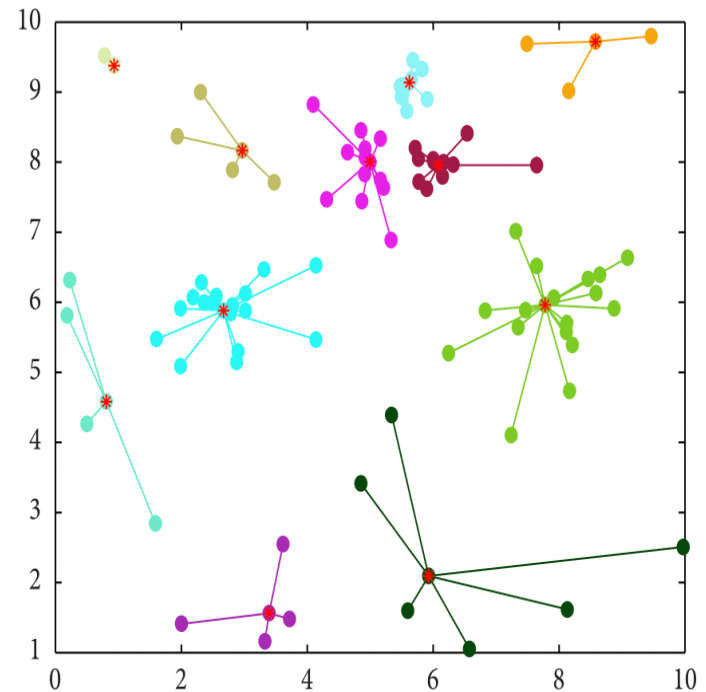
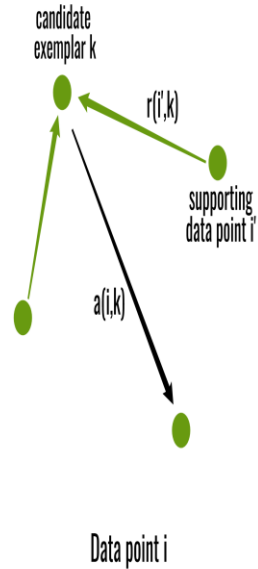
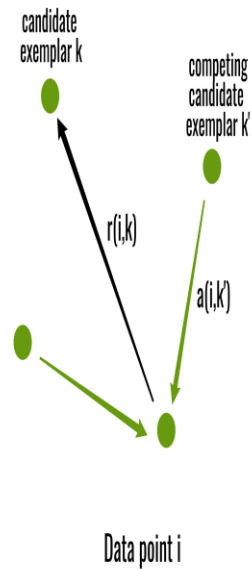
Clustering Algorithms: Detailed Explanation

Affinity Propagation, Mean Shift, Spectral,
DBSCAN, OPTICS, and Birch

Affinity Propagation Clustering

- Affinity propagation is based on 'message passing' between data points to find exemplars and form clusters. The algorithm identifies a high-quality set of exemplars and corresponding clusters without pre-specifying the number of clusters.
- Pros: Auto-determines cluster count, handles non-convex shapes.
- Cons: Computationally intensive for large datasets.


```
from sklearn.cluster import AffinityPropagation
aff = AffinityPropagation(random_state=42)
y_aff = aff.fit_predict(X)
```



Mean Shift Clustering

- Mean shift identifies 'blobs' in data by shifting data points towards areas of higher density iteratively. The bandwidth parameter impacts the outcome significantly.
- Pros: No need to specify cluster count, works for non-linear clusters.
- Cons: Bandwidth selection is critical, expensive for high-dimensional data.

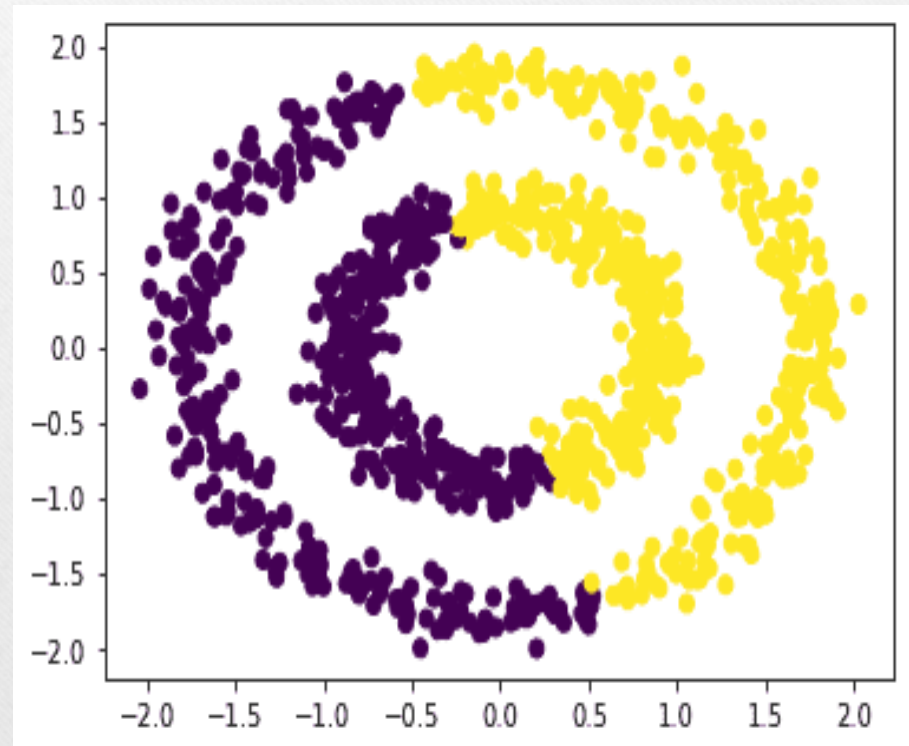
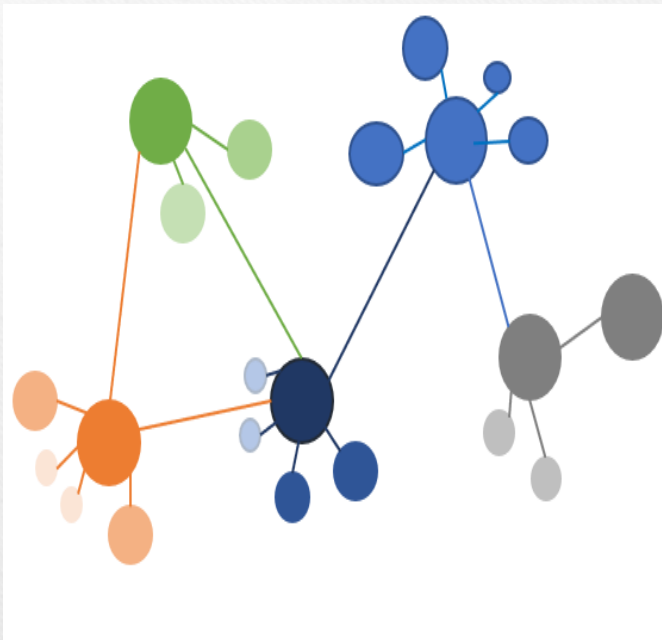

```
from sklearn.cluster import MeanShift
mean_shift = MeanShift()
y_mean_shift = mean_shift.fit_predict(X)
```



Spectral Clustering

- Spectral clustering uses the eigenvalues of a similarity matrix for dimensionality reduction and clustering in a lower-dimensional space. It is suitable for non-convex clusters.
- Pros: Works well with complex shapes.
- Cons: Sensitive to the similarity matrix and the number of clusters.

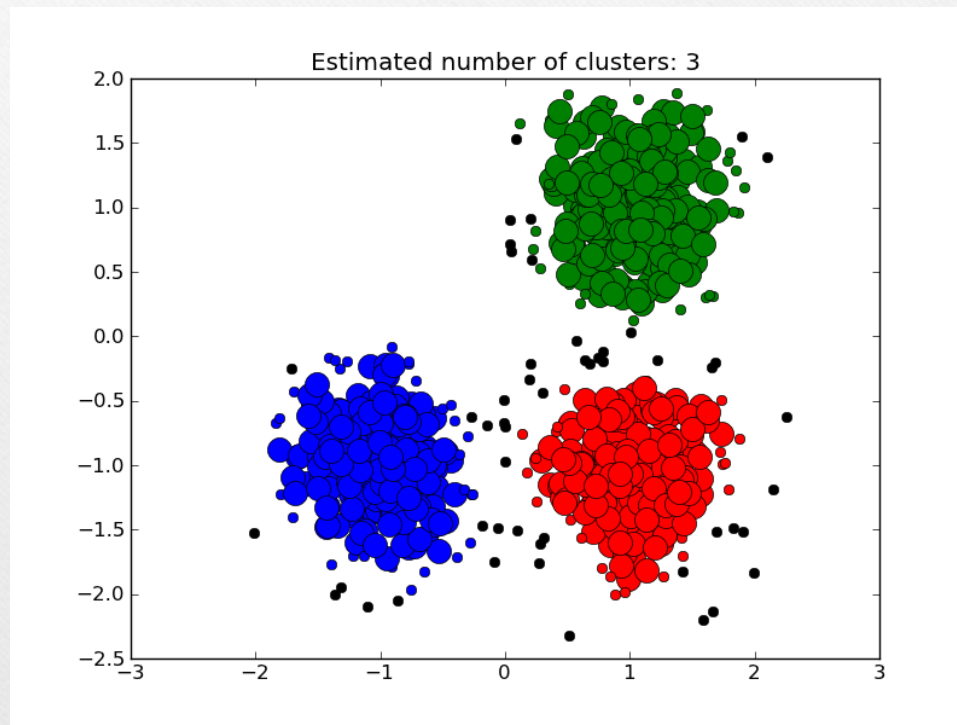

```
from sklearn.cluster import SpectralClustering  
spectral = SpectralClustering(n_clusters=3,  
random_state=42)  
y_spectral = spectral.fit_predict(X)
```



DBSCAN Clustering

- DBSCAN clusters based on a density criterion and distance metric, identifying outliers as noise. It uses two parameters: 'eps' and 'minPts'.
- Pros: Detects clusters of varying shapes, identifies noise.
- Cons: Depends on parameter choices, not ideal for varying-density datasets.

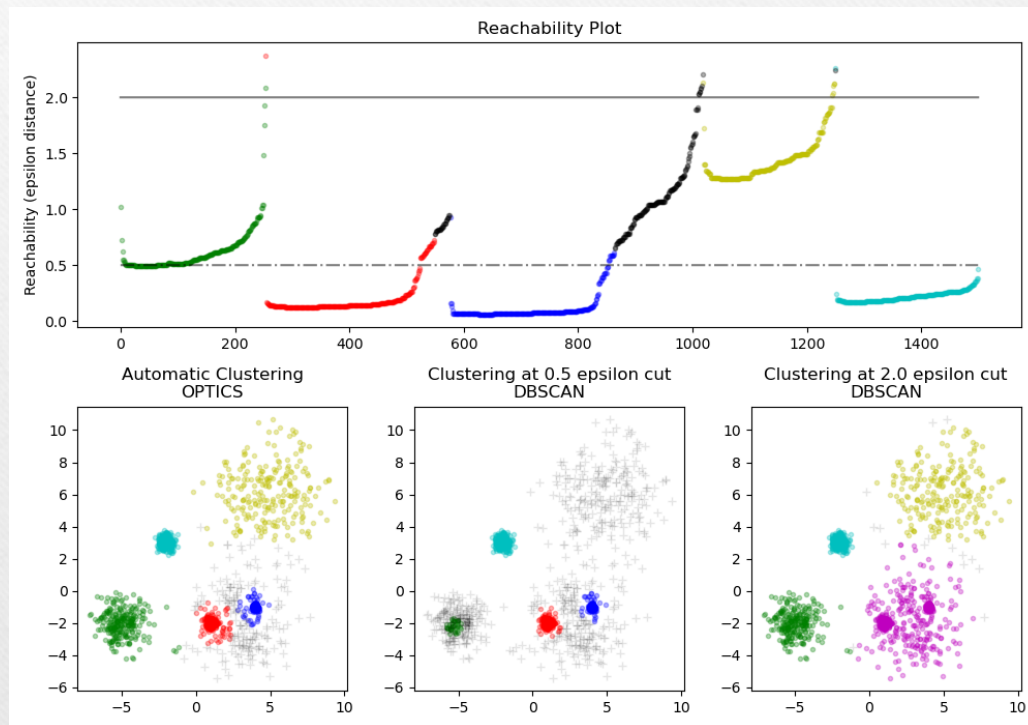

```
from sklearn.cluster import DBSCAN  
dbscan = DBSCAN(eps=0.5, min_samples=5)  
y_dbscan = dbscan.fit_predict(X)
```



OPTICS Clustering

- OPTICS handles varying-density clusters by ordering data points to create a reachability plot. This plot helps identify clusters at different density levels.
- Pros: Detects clusters with varying density, hierarchical structure.
- Cons: Requires post-processing to extract clusters.


```
from sklearn.cluster import OPTICS
opt= OPTICS(min_samples=5,max_eps=np.inf,metric='minkowski',
p=2)
y_optics = optics.fit_predict(X)
```

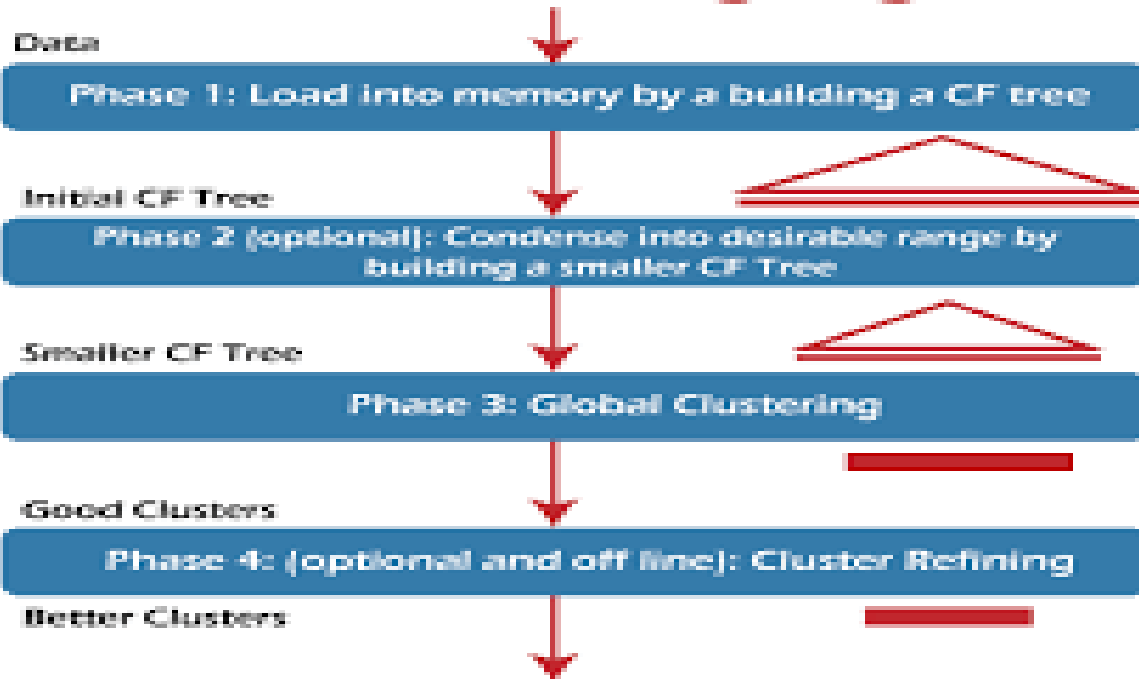


Birch Clustering

- Birch builds clusters using a tree structure (CF tree), suitable for large datasets and incremental clustering.
- Pros: Scales well, efficient for large data.
- Cons: Sensitive to input order, depends on branching factor and threshold.


```
from sklearn.cluster import Birch
birch = Birch(threshold=0.5, n_clusters=None)
y_birch = birch.fit_predict(X)
```

The BIRCH Clustering Algorithm



Birch Clustering

