

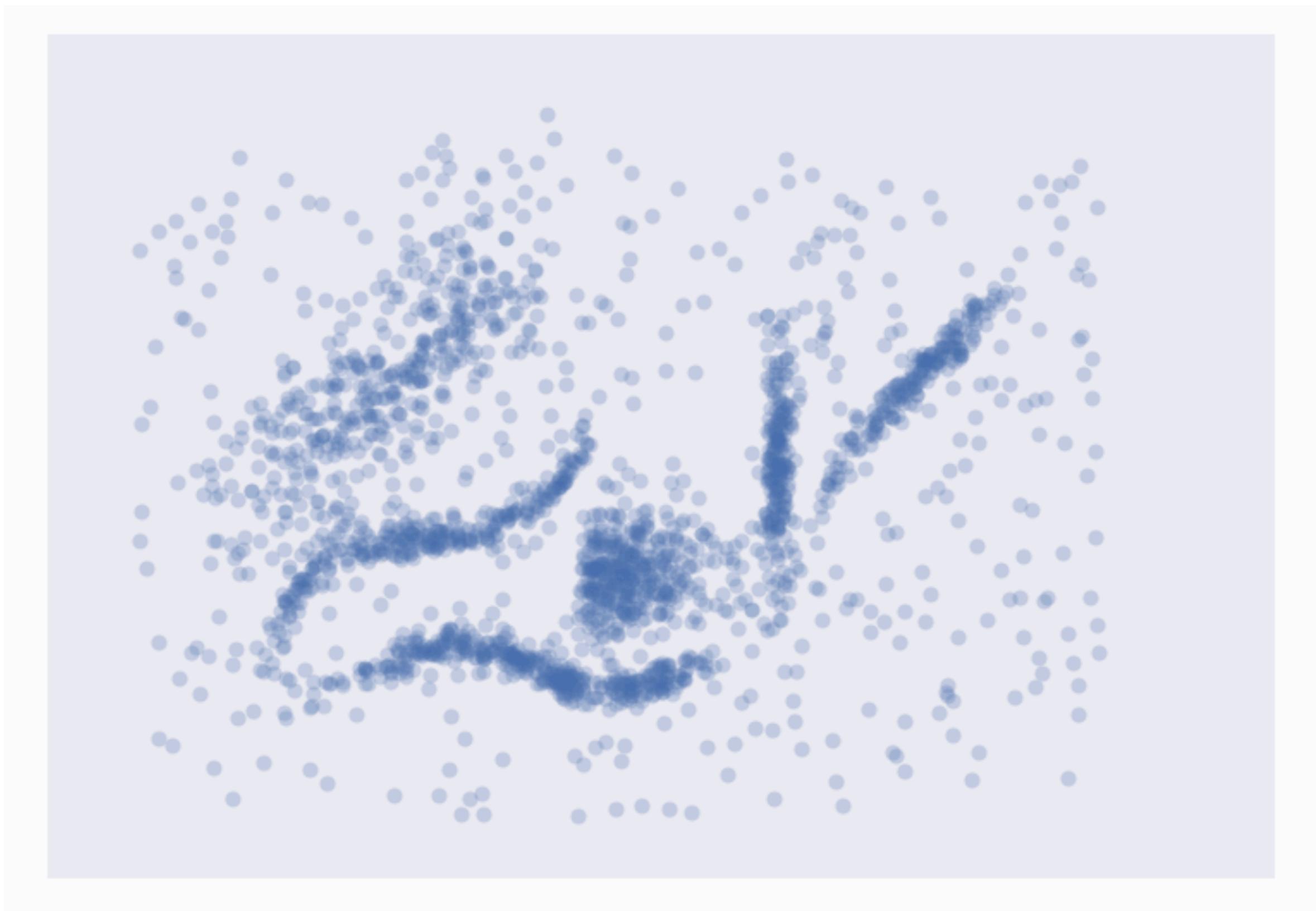
DBSCAN – Density-based spatial clustering of applications with noise

HDBSCAN – Hierarchical DBSCAN

Lu Li
20210416

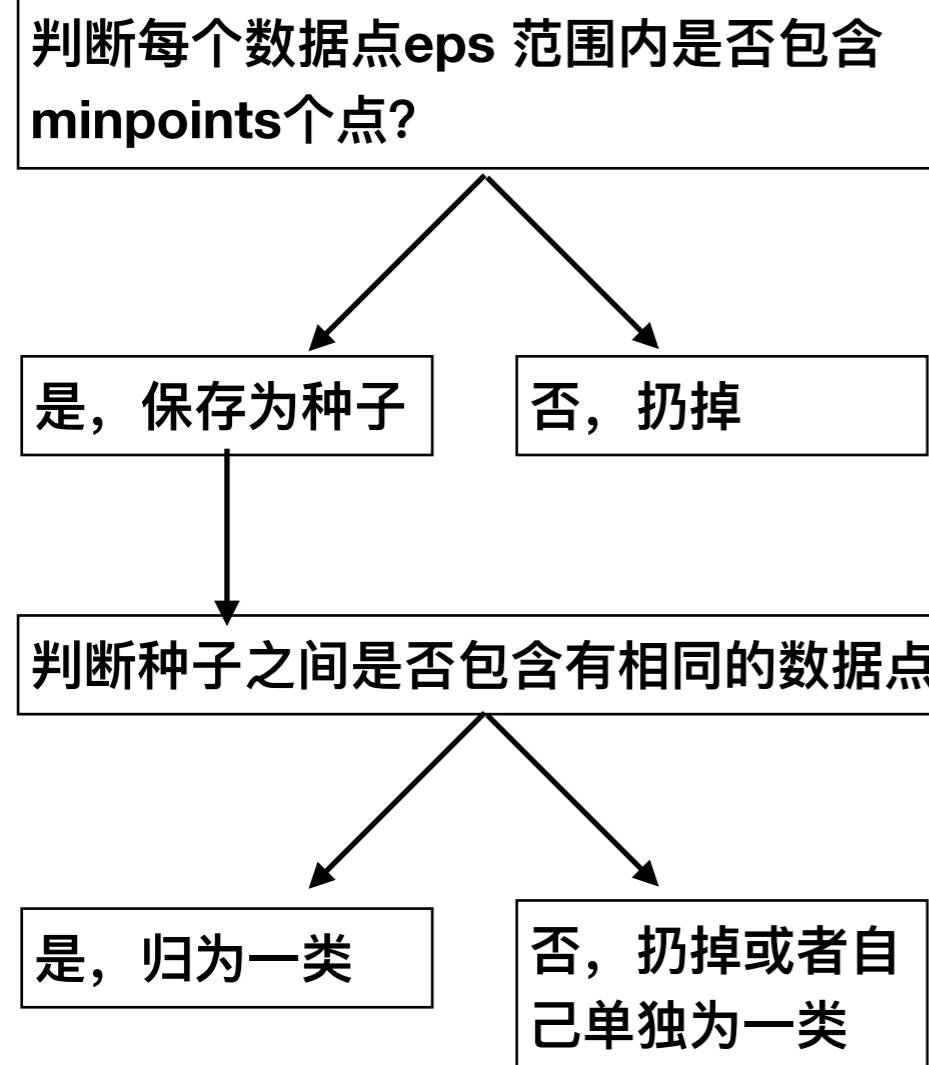
https://hdbSCAN.readthedocs.io/en/latest/how_hdbSCAN_works.html

Clustering

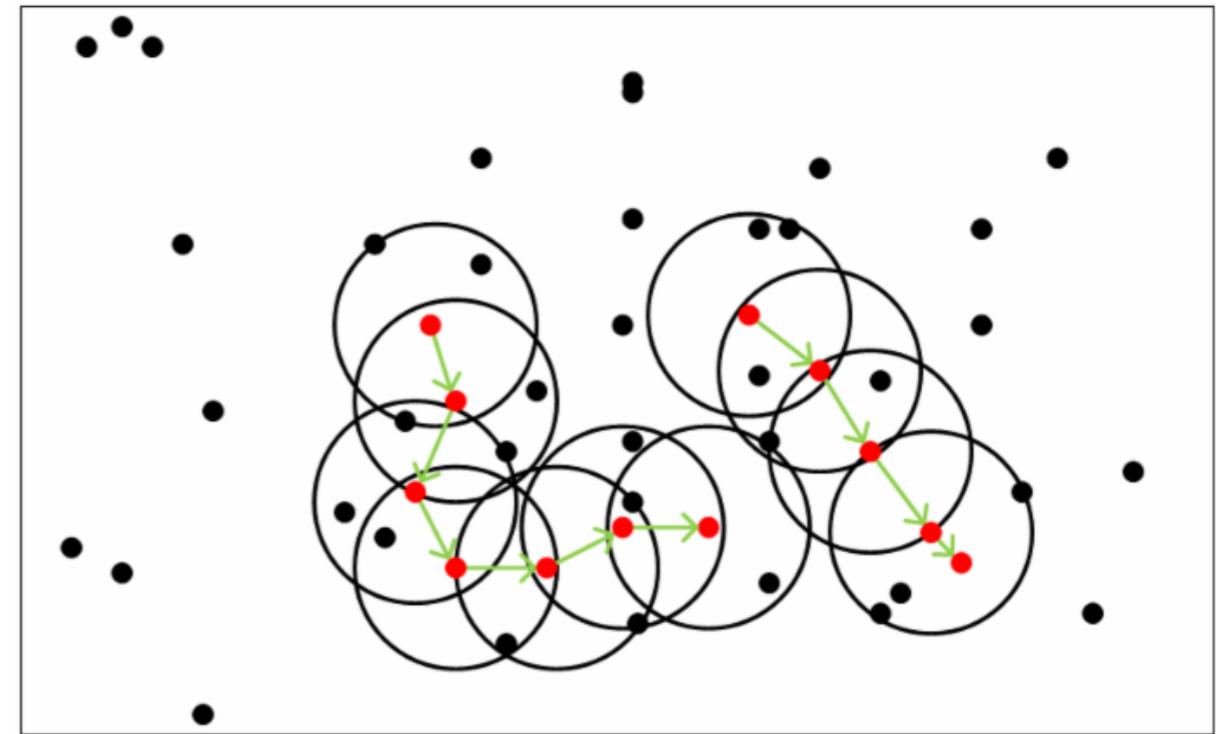


DBSCAN

参数: eps (邻域的半径)
 minPoints (邻域半径内包含的点数)



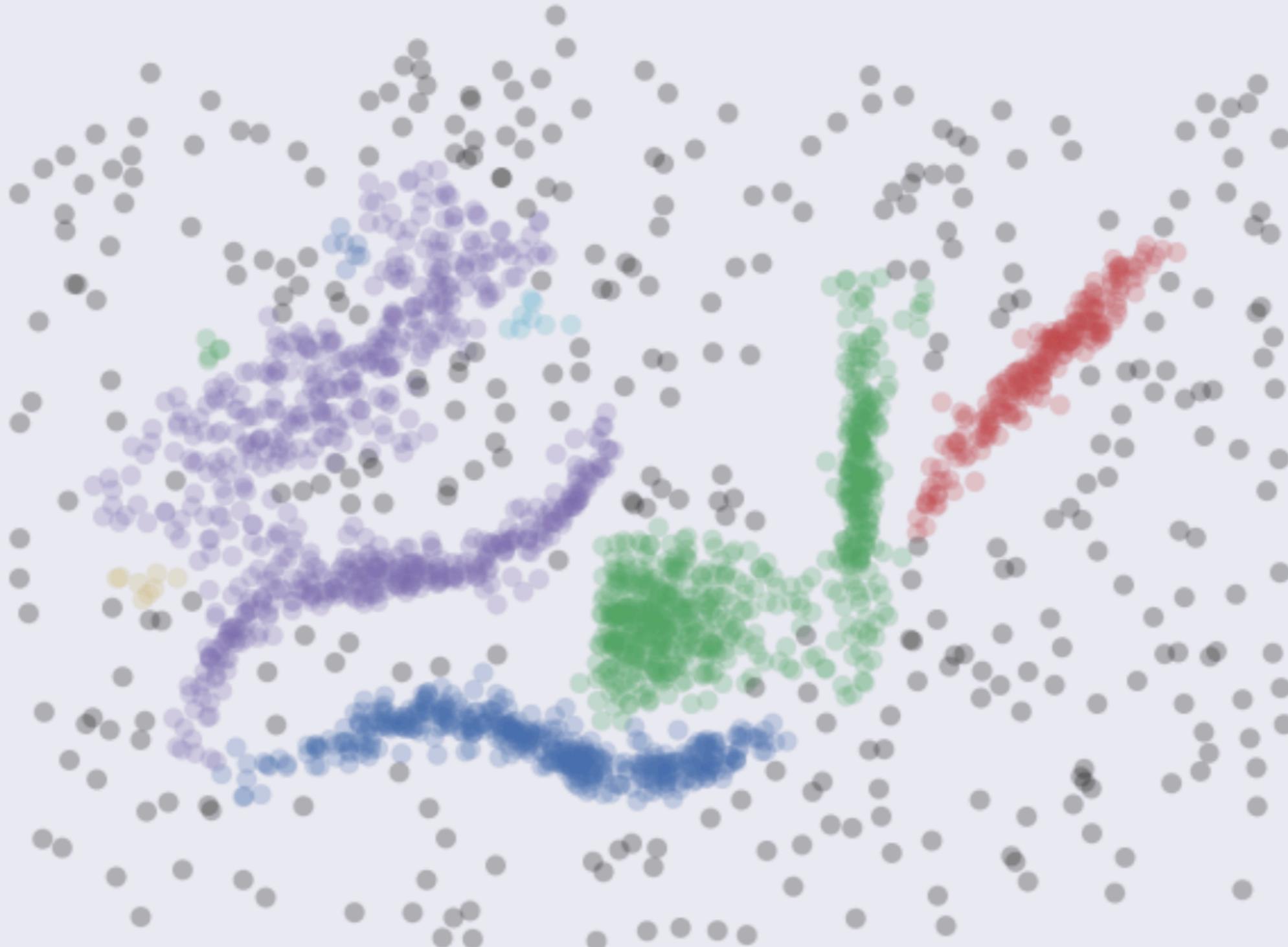
FOF: $\text{minPoints} = 1$



定义了临界密度一刀切

Clusters found by DBSCAN

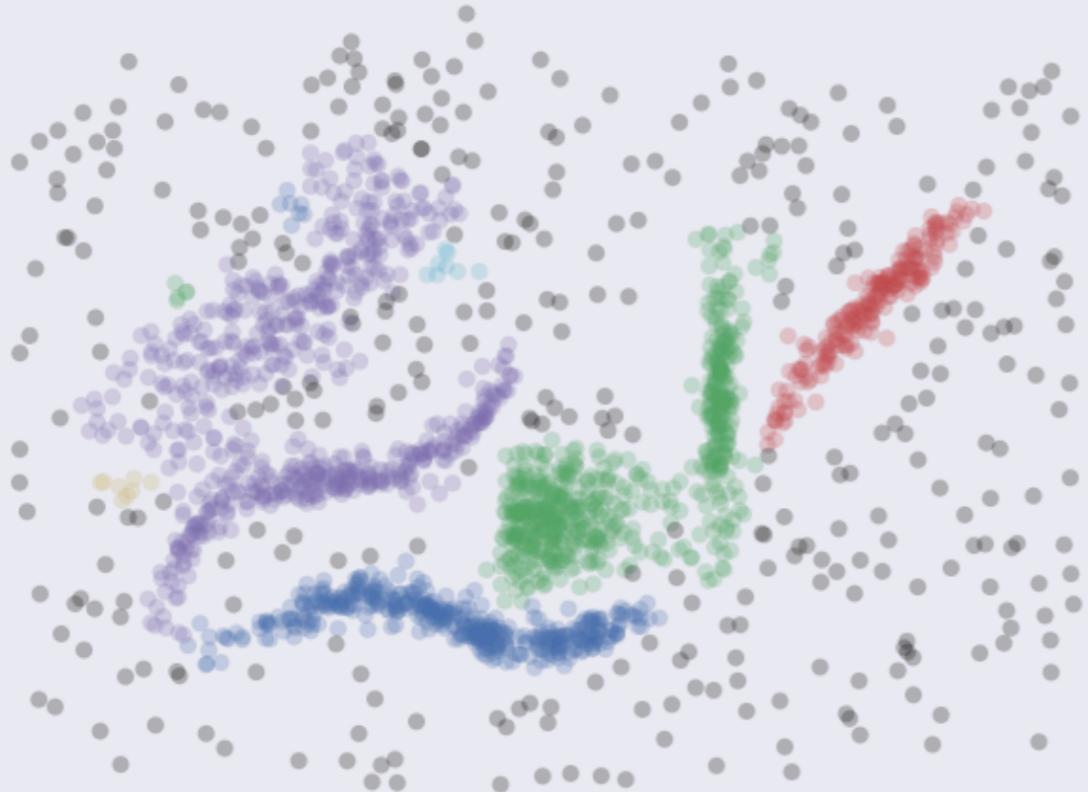
Clustering took 0.02 s



HDBSCAN

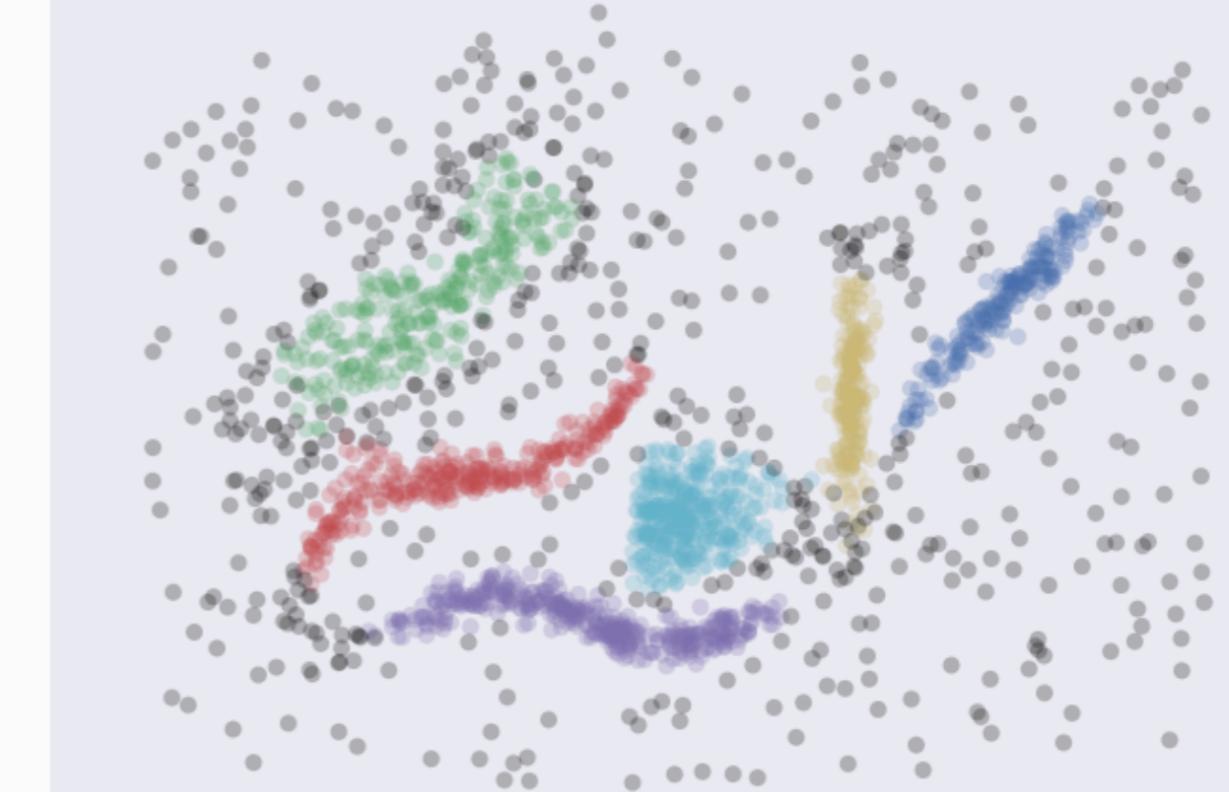
Clusters found by DBSCAN

Clustering took 0.02 s



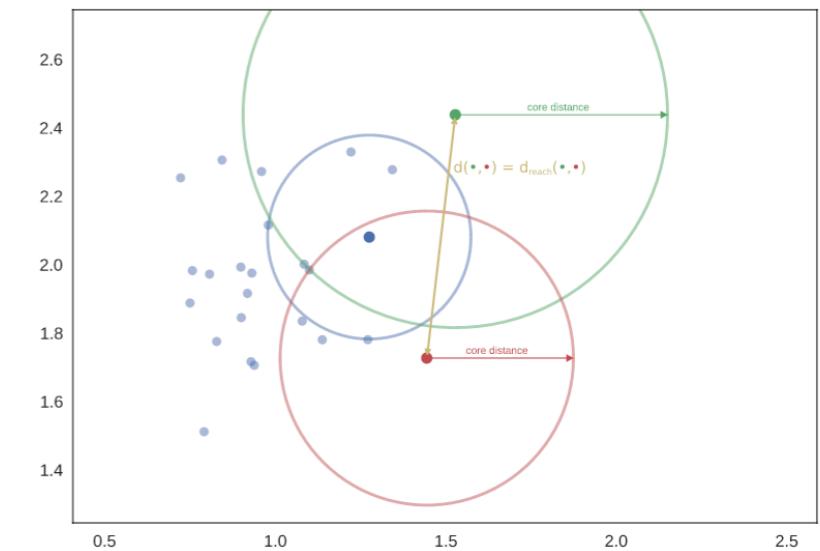
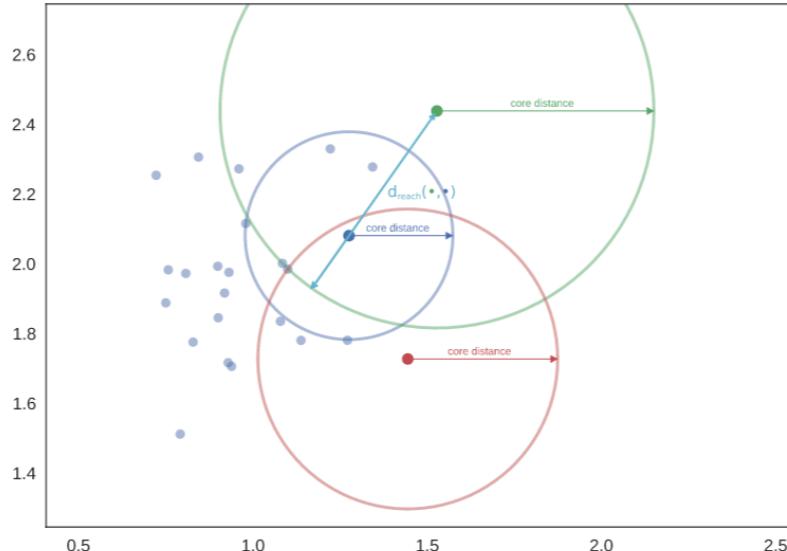
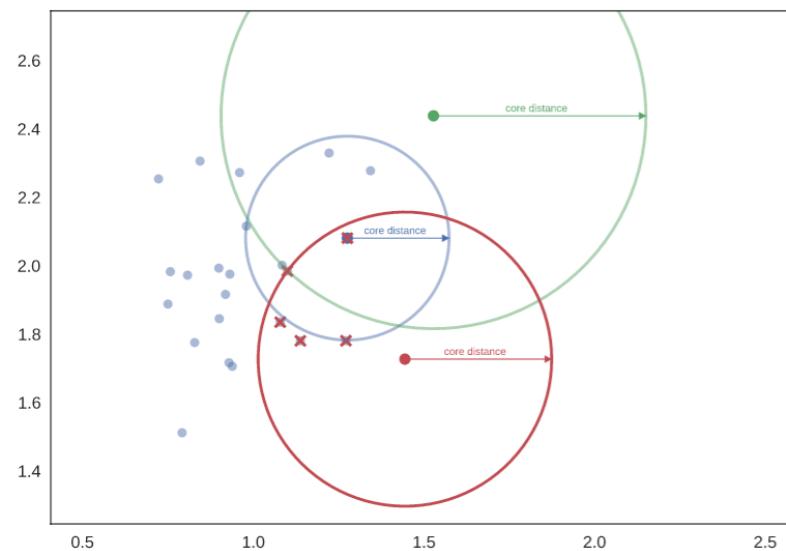
Clusters found by HDBSCAN

Clustering took 0.06 s



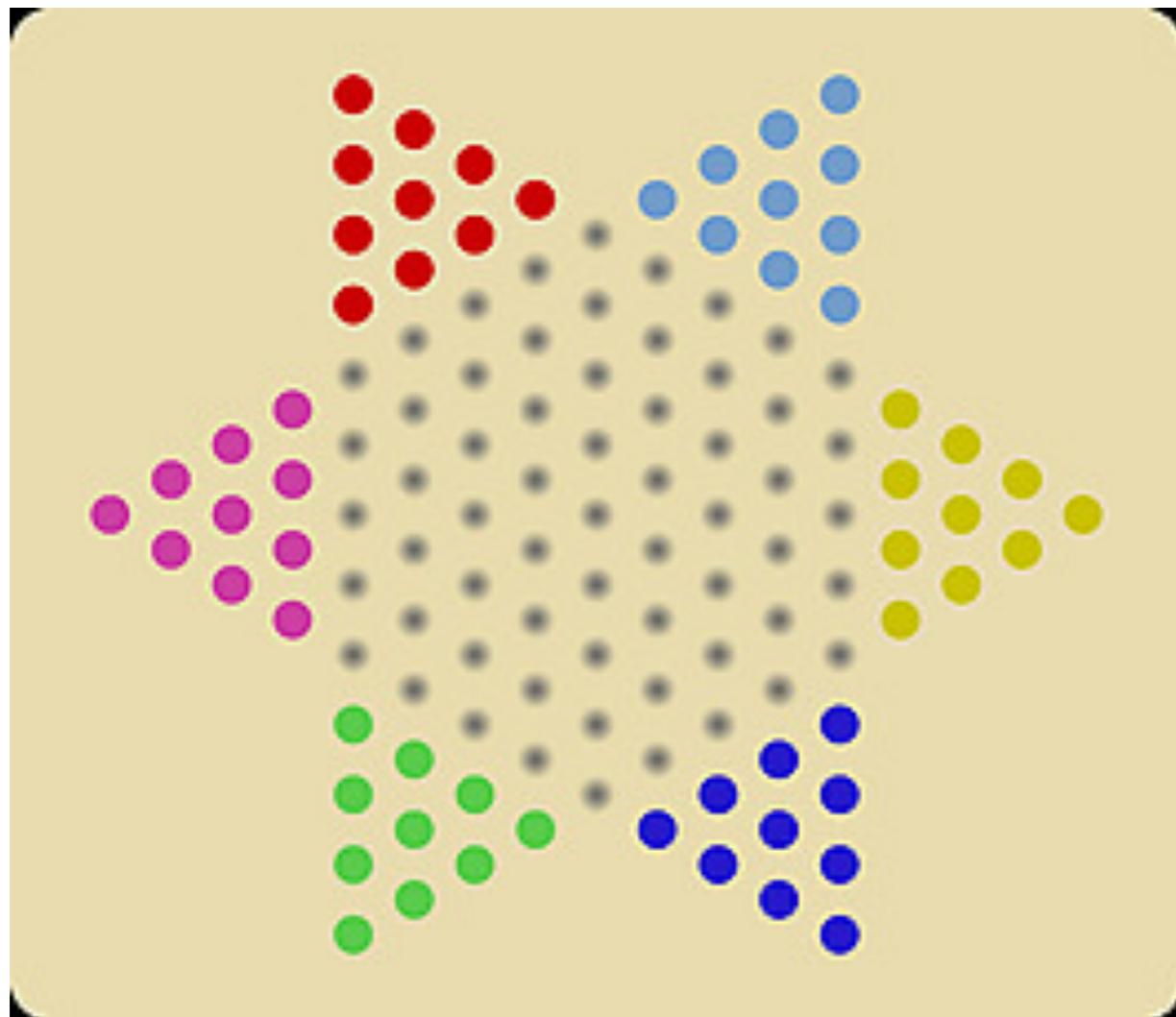
mutual reachability distance

$$d_{mreach-k}(a, b) = \max \{ \text{core}_k(a), \text{core}_k(b), d(a, b) \}$$



mutual reachability distance

$$d_{mreach-k}(a, b) = \max \{ \text{core}_k(a), \text{core}_k(b), d(a, b) \}$$

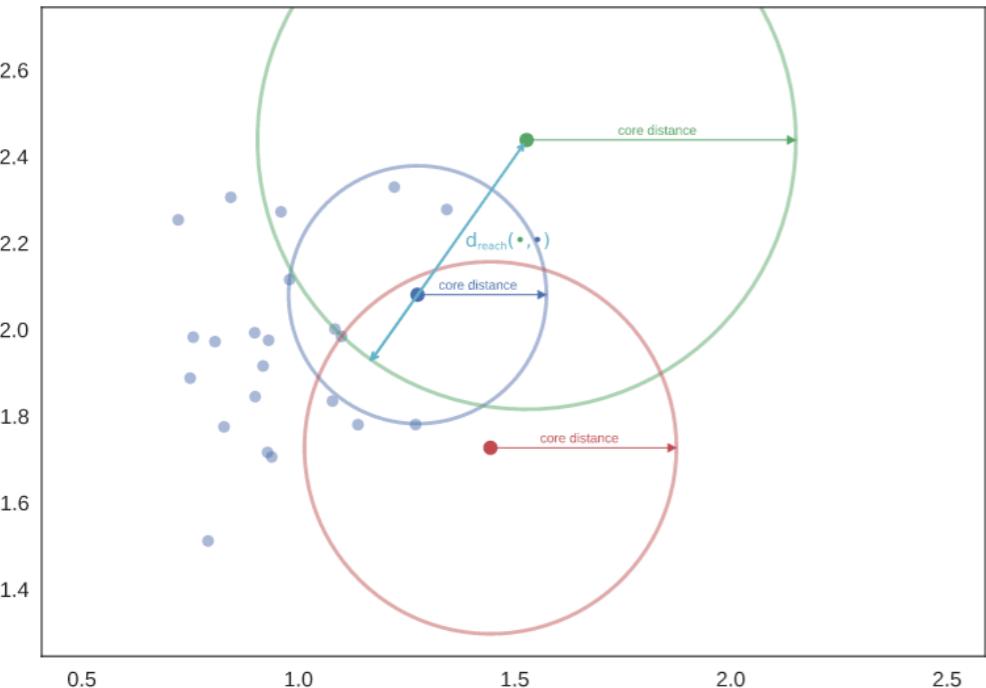


空间密度均匀：

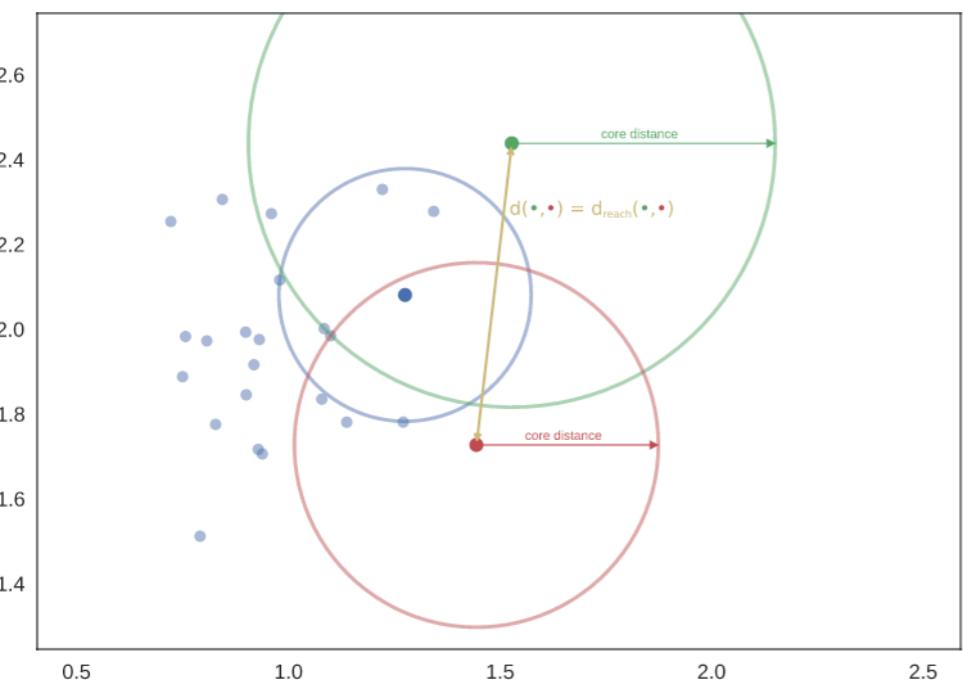
点与点之间互相可达距离一致

mutual reachability distance

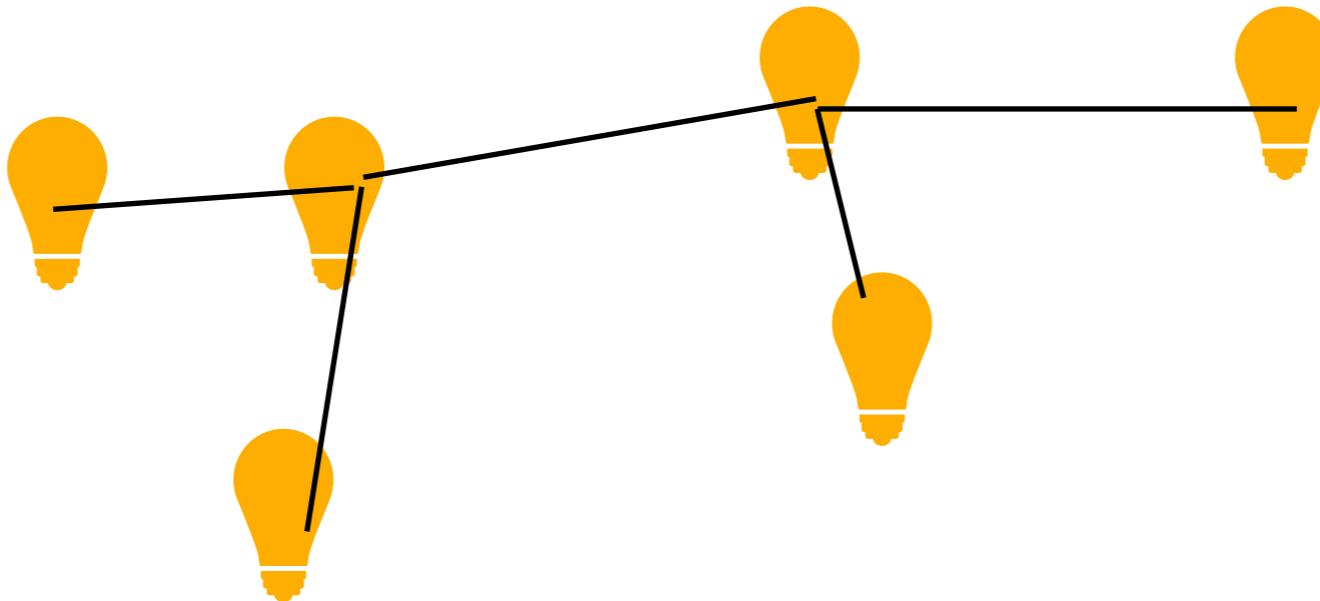
$$d_{mreach-k}(a, b) = \max \{ \text{core}_k(a), \text{core}_k(b), d(a, b) \}$$



同一个簇内的点，相互距离没有影响；
噪声与噪声点之间的距离也没有影响。
但是簇和簇之间的噪声点被“推开”。



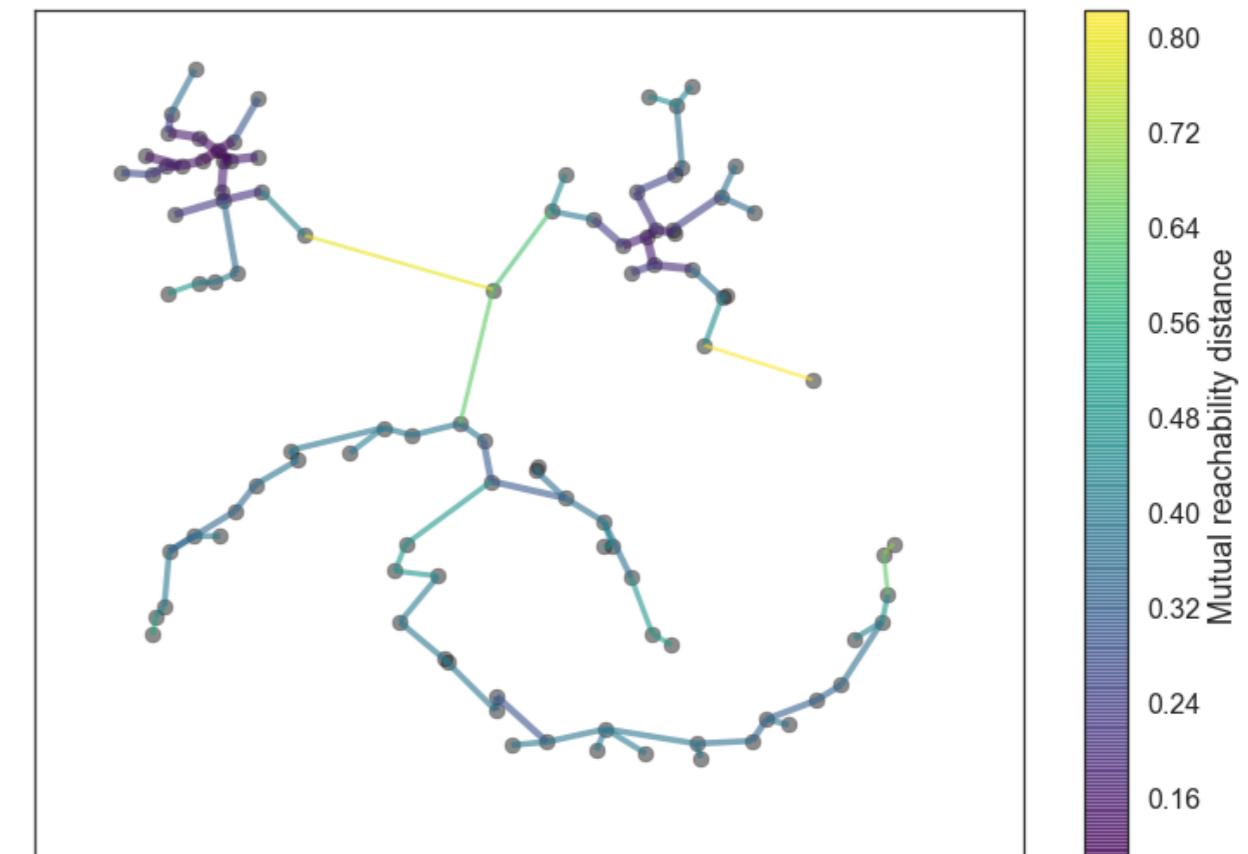
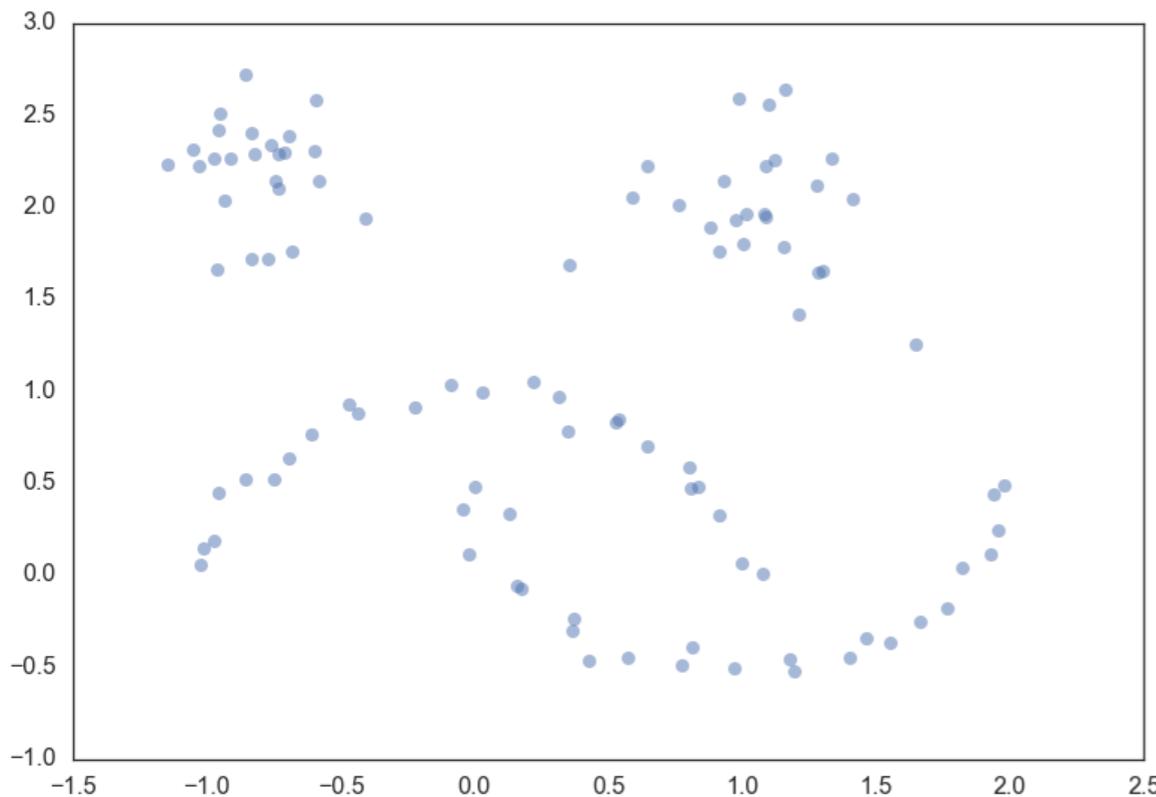
Minimum Spanning Tree



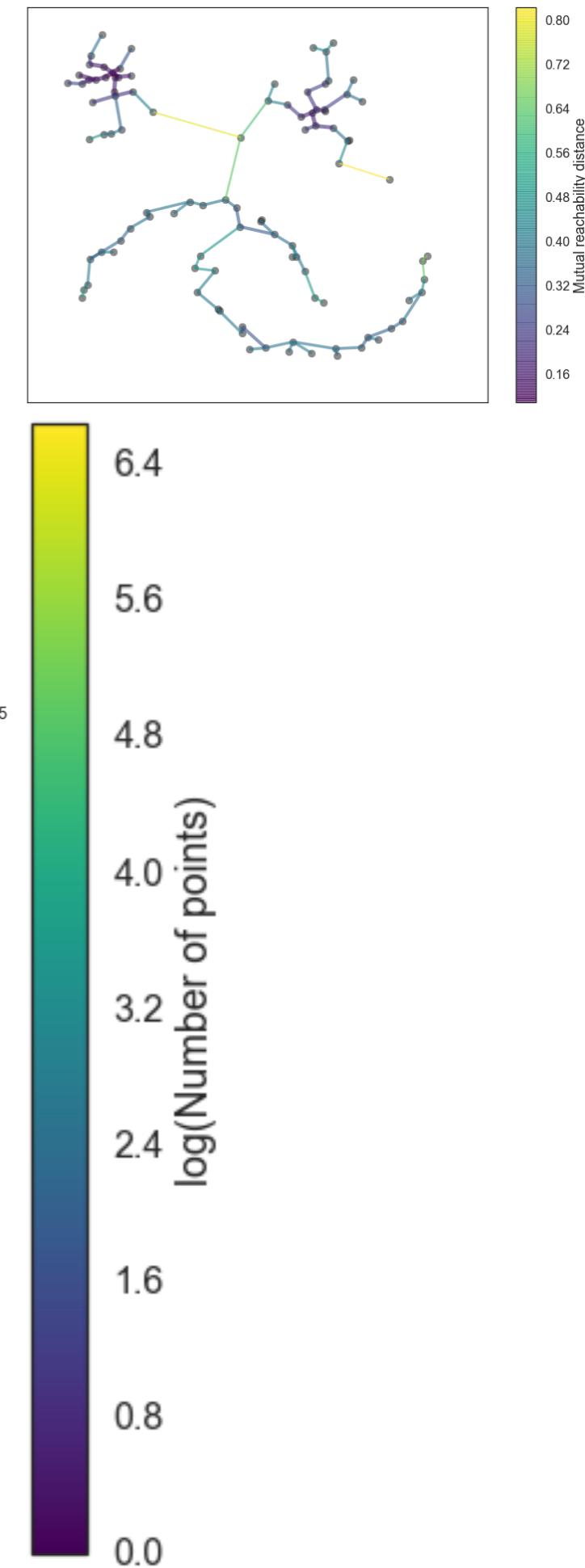
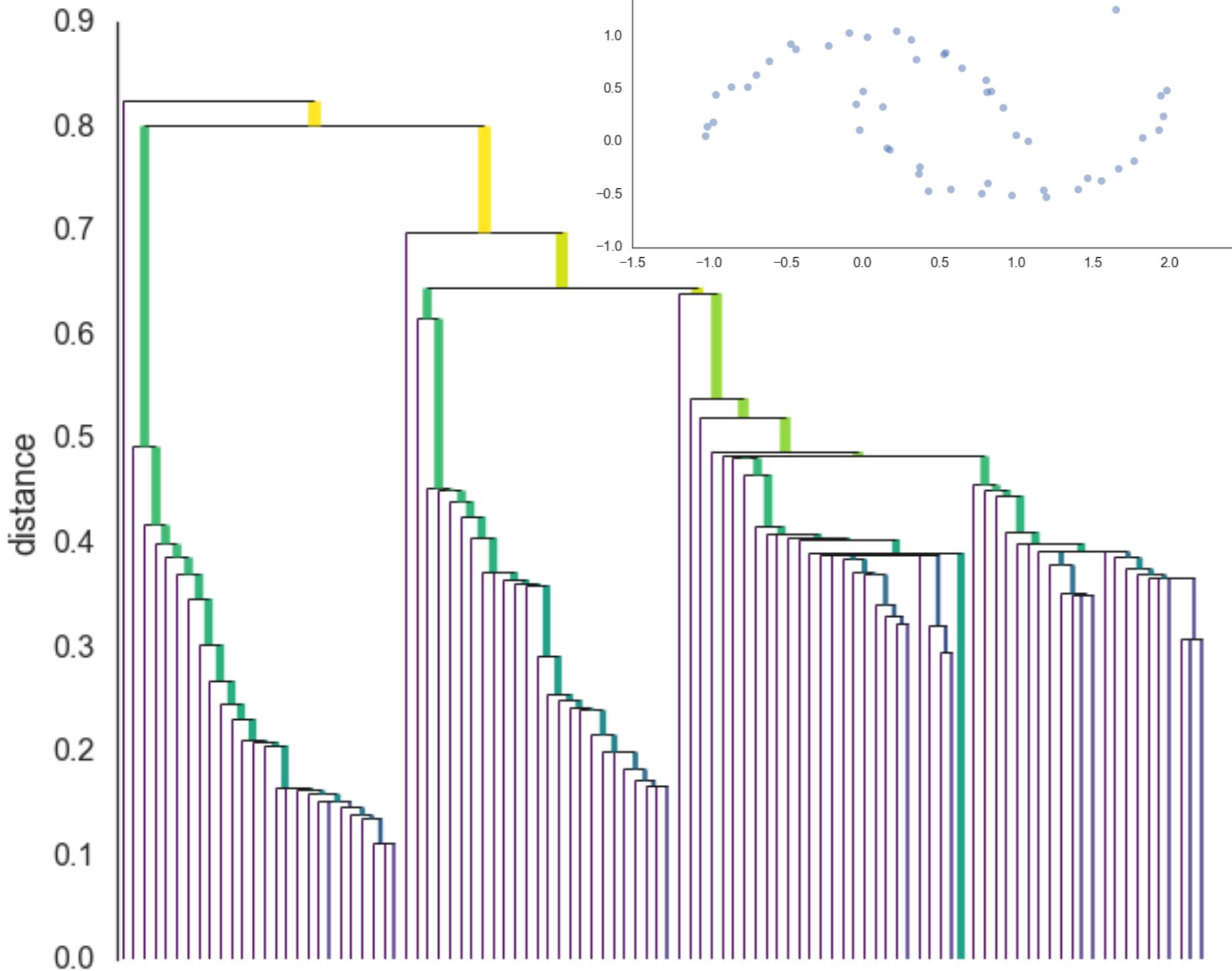
最小生成树

要在 n 个城市之间铺设光缆，主要目标是要使这 n 个城市的任意两个之间都可以通信，使铺设光缆的总费用最低。

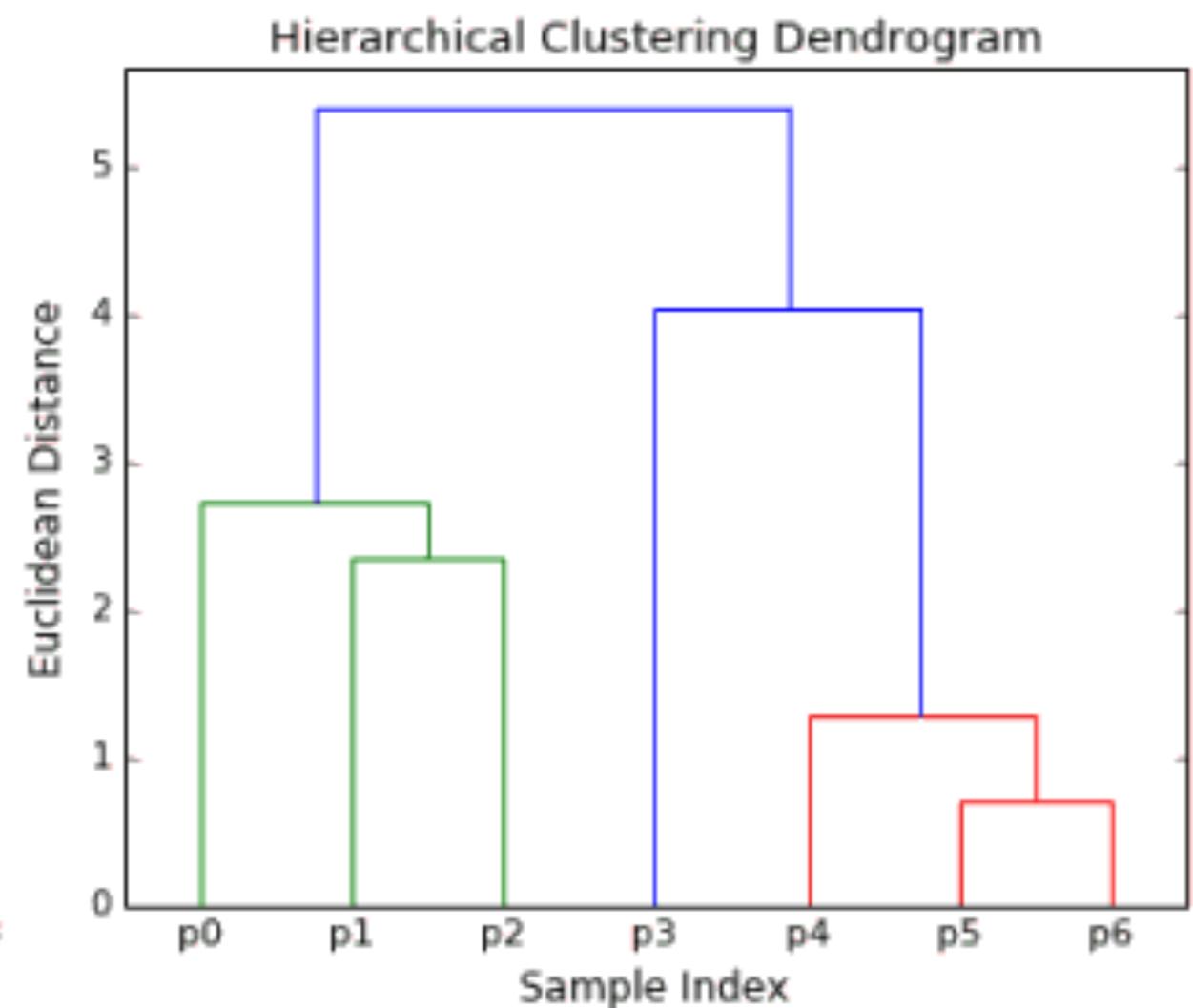
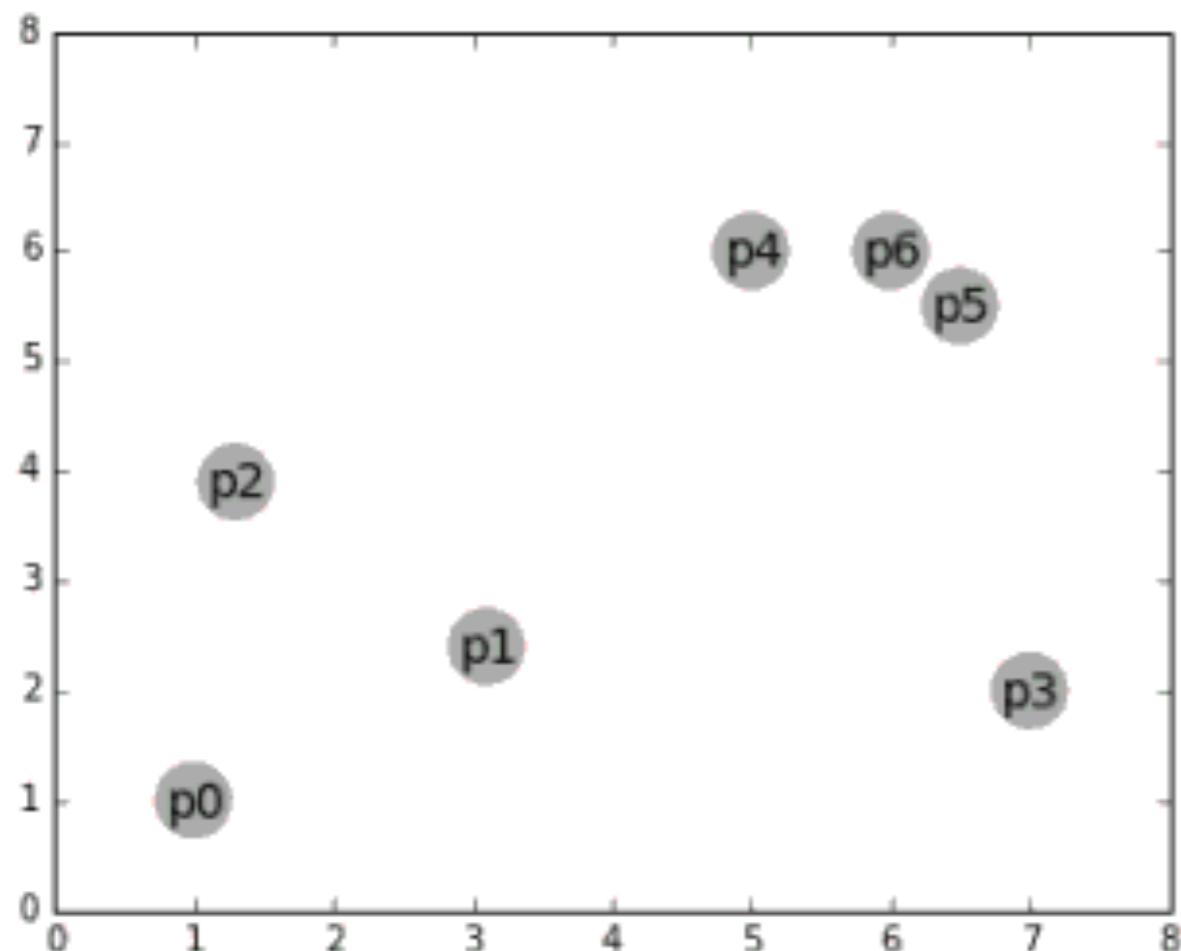
Minimum Spanning Tree



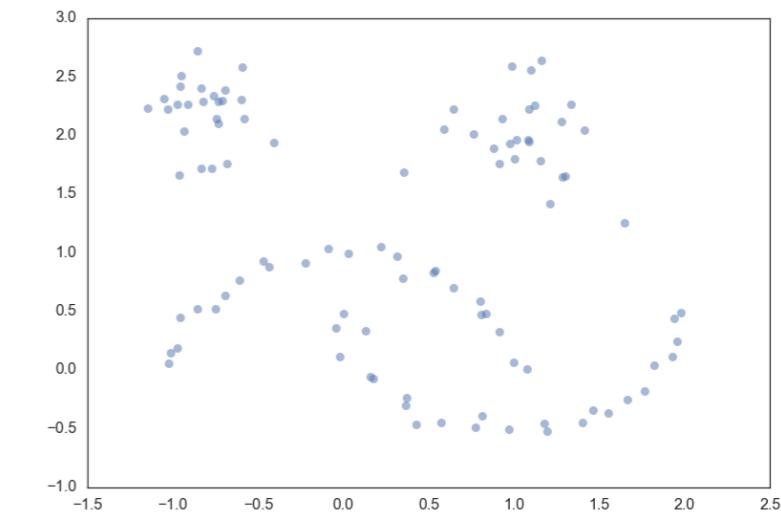
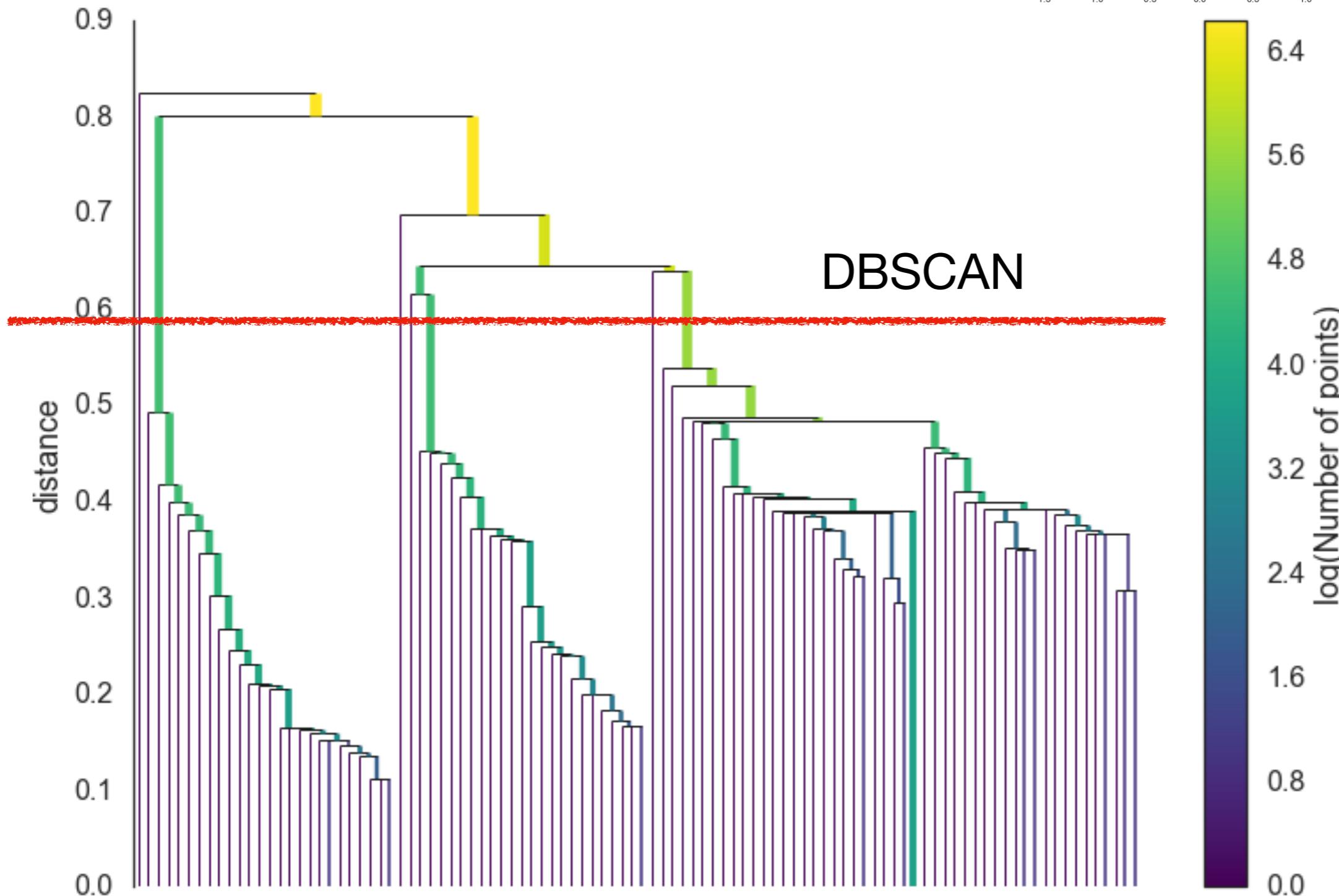
Hierarchical Clustering



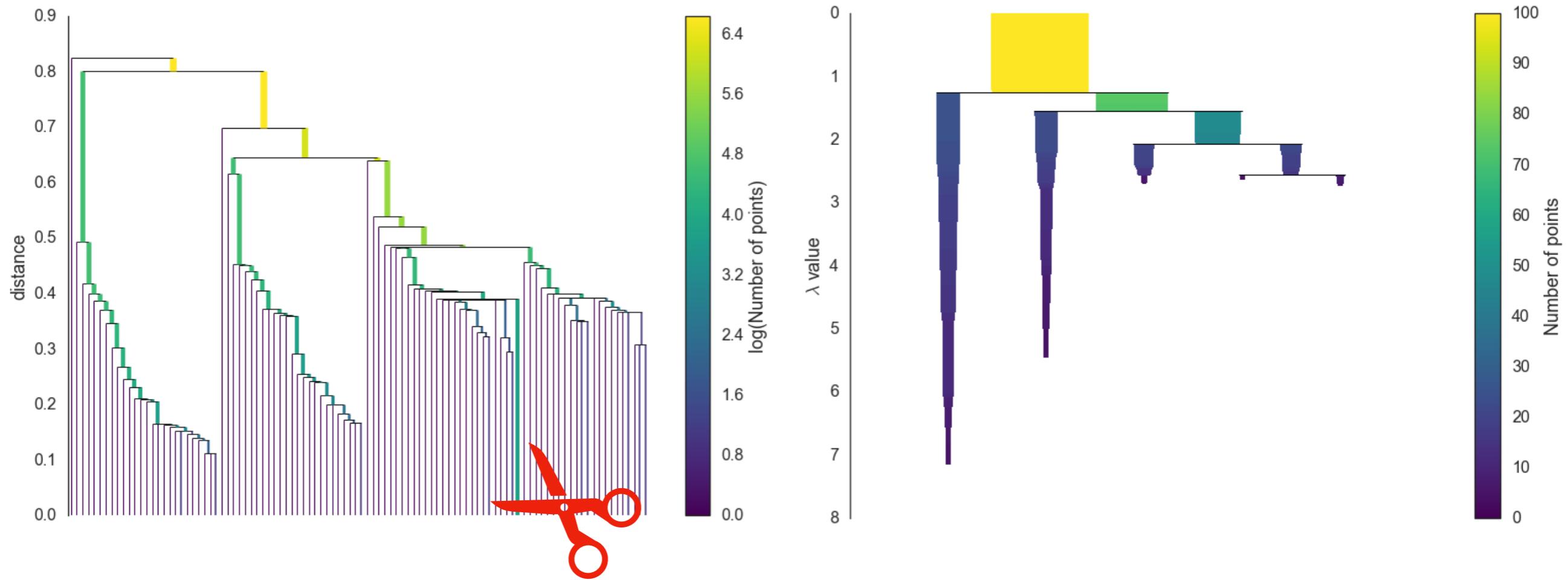
Hierarchical Clustering



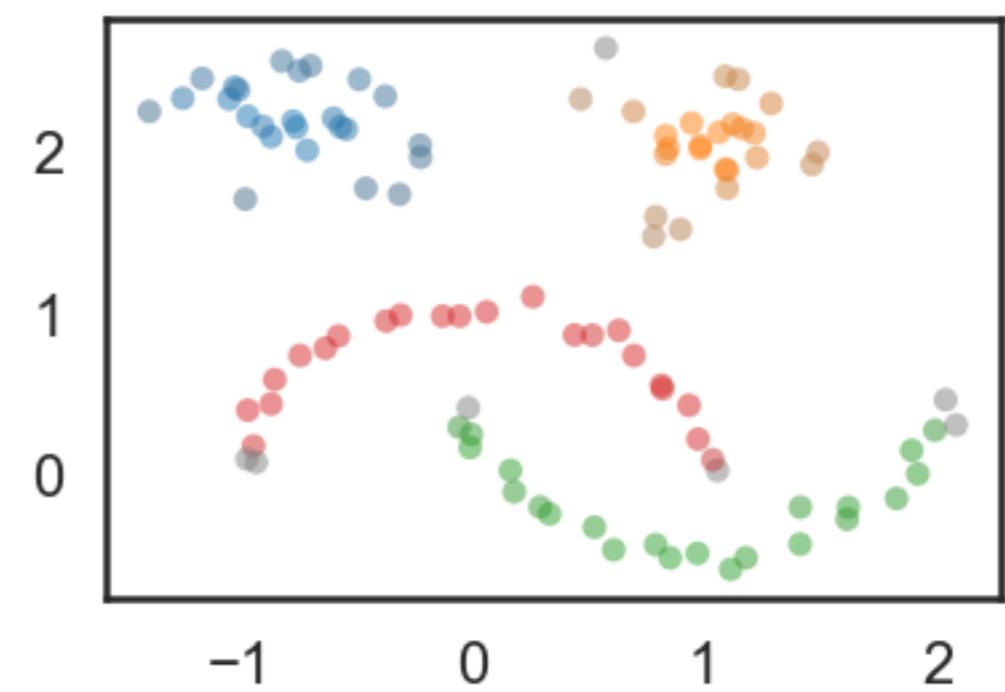
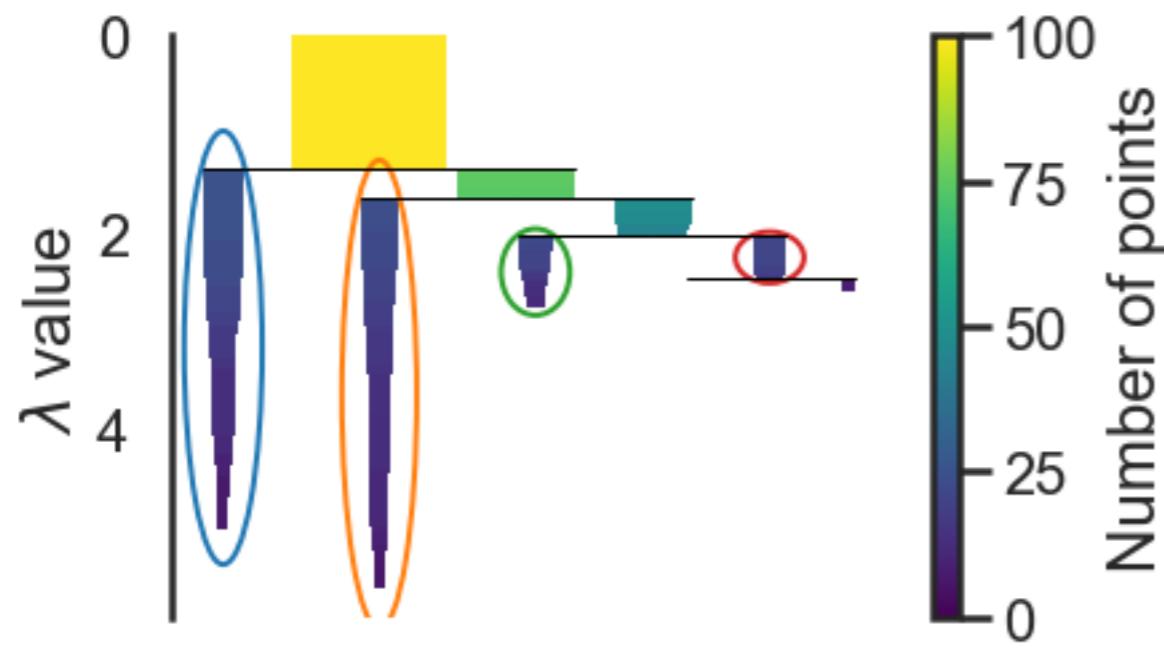
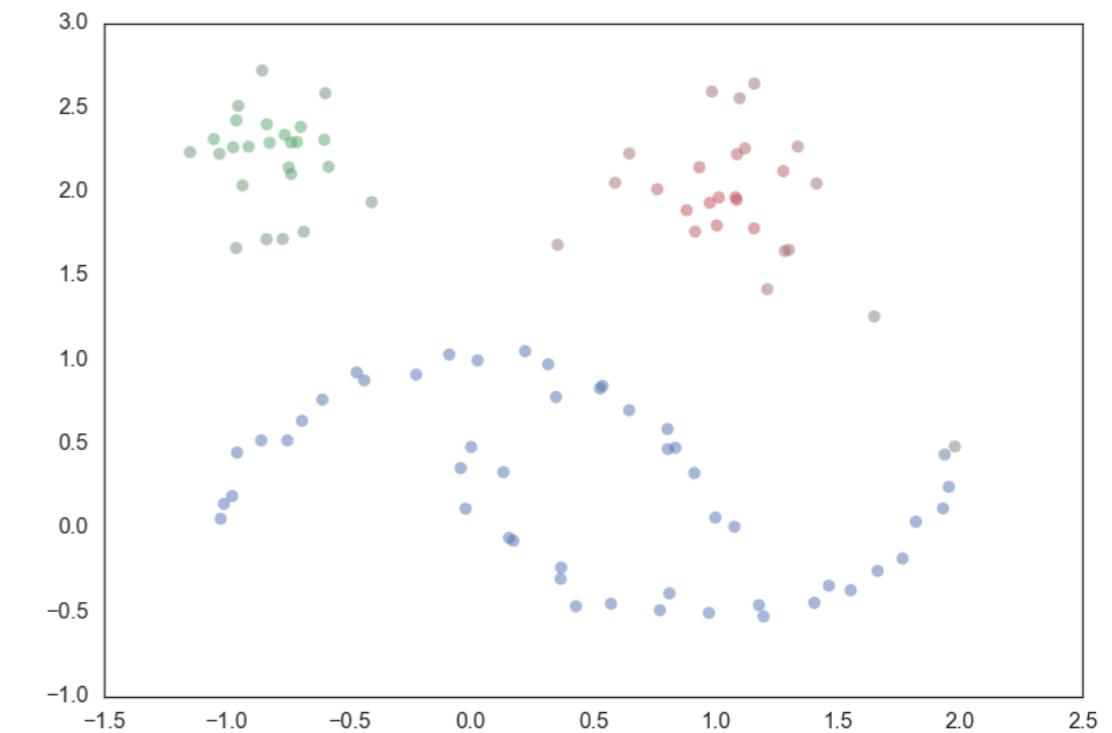
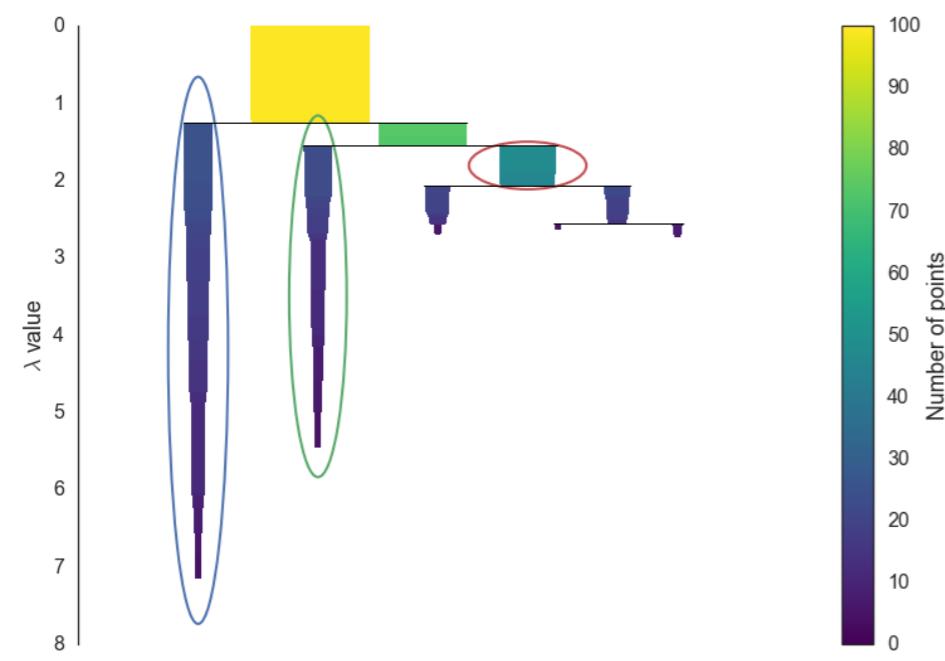
Hierarchical Clustering



Condense the cluster tree

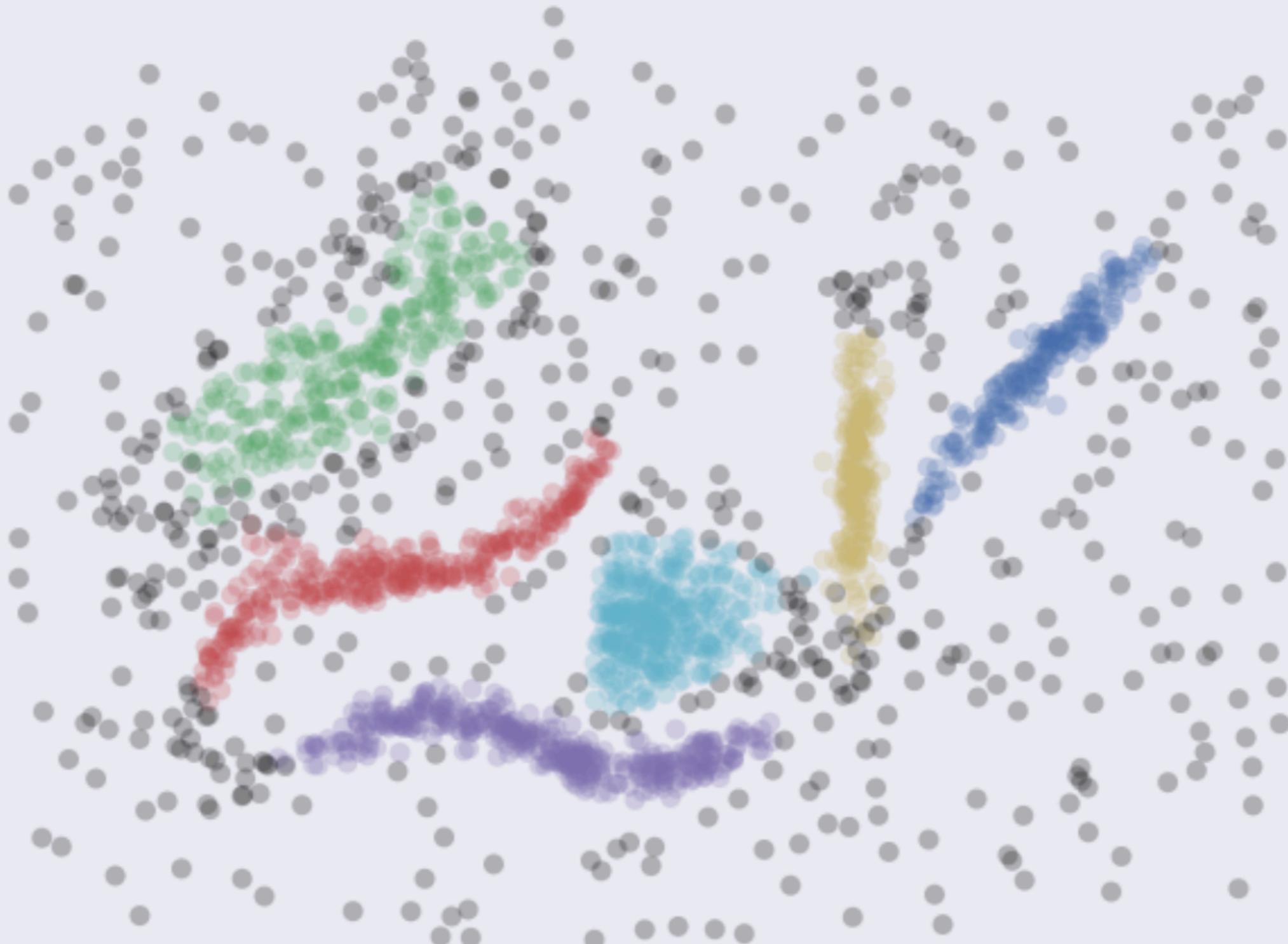


修剪枝丫：从上至下，历遍节点，剪掉小于 min cluster members 的节点



Clusters found by HDBSCAN

Clustering took 0.06 s



使用相互可达距离替换欧氏距离，该距离可以使得密度低的点离密度高的区域更远，减少
dbSCAN对Eps阈值的依赖性

使用最小生成树构建层次聚类模型，引入层次聚类思想

对最小生成树的最小子树做了限制，减少计算量，同时保证生成的类簇不要过小

使用“簇稳定性”的度量方式自动划分类簇，不需要自行设定阈值

