

## 2. Ý nghĩa tham số radius, min sample trong thuật toán dbscan? Nếu chỉ số lớn, nhỏ ảnh hưởng thế nào tới thuật toán?

Radius: Một giá trị khoảng cách được sử dụng để xác định *vùng lân cận radius* của bất kỳ điểm dữ liệu nào.

Nếu Radius được chọn quá nhỏ, một phần lớn dữ liệu sẽ không được phân cụm và được xem là *noise*; trong khi đối với giá trị Radius quá cao, các cụm sẽ hợp nhất và phần lớn các điểm sẽ nằm trong cùng một cụm.

Min sample: Là một ngưỡng số điểm dữ liệu tối thiểu được nhóm lại với nhau nhằm xác định một *vùng lân cận radius* có mật độ cao. Số lượng Min sample không bao gồm điểm ở tâm.

Nếu Min sample thấp có nghĩa là nó sẽ xây dựng nhiều cụm hơn từ là tiếng ồn( Min sample ít nhất là 3), trong khi Min sample cao có nghĩa nó có thể bao gồm tất cả các điểm trong tập dữ liệu( nên chọn Min sample =  $2 * \text{dim}$ , dim là chiều của dữ liệu, còn đối với dữ liệu 2 chiều nên chọn Min sample = 4).

## 3. So sánh ba thuật toán: kmean, GMM, dbscan. Khi nào nên sử dụng thuật toán nào? cho ví dụ?

Model	Pros	Cons	Use Cases	Example
K-mean	<ul style="list-style-type: none"><li>+) Quickest centroid based algorithm</li><li>+) Very lucid and can scale up for large amount of dataset</li><li>+) Reduces intra-cluster variance measure</li></ul>	<ul style="list-style-type: none"><li>+) Suffers when there is noise in the data</li><li>+) Outliers can never be identified</li><li>+) Even though it reduces intra-cluster variance, it faces local minimum problem</li><li>+) Not ideal for datasets of non-convex shapes</li><li>+) Complicated to predict best K value</li></ul>	Even cluster size, flat geometry, not too many clusters and general-purpose	<ul style="list-style-type: none"><li>+) Segmentation of customers in business</li><li>+) Image segmentation</li><li>+) Genetic analysis in medicine</li><li>+) Anomaly detection.</li></ul>
DBSCAN	<ul style="list-style-type: none"><li>+) Resistant to outliers</li><li>+) Can handle clusters of different shapes and sizes</li><li>+) Not required to specify the number of cluster</li></ul>	<ul style="list-style-type: none"><li>+) Highly sensitive to the two-parameters: Radius and Min sample</li><li>+) DBSCAN cannot cluster datasets well with large variances in densities</li></ul>	Uneven cluster sizes and non-flat geometry	<ul style="list-style-type: none"><li>+) Face clustering</li></ul>
GMM	<ul style="list-style-type: none"><li>+) Robust to outliers</li><li>+) Provides the BIC score for selecting parameters</li><li>+) Converges fast given good initialization</li></ul>	The algorithm is highly complex and can be slow	Good for density estimation and flat geometry	<ul style="list-style-type: none"><li>+) Clustering of homogeneous bacterial colonies to estimate their size</li></ul>